

**PCT**WORLD INTELLECTUAL PROPERTY ORGANIZATION  
International Bureau

## INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

<b>(51) International Patent Classification <sup>6</sup> :</b> <b>C12N 15/11, 15/10, C07K 14/47, C12P 21/00, C12Q 1/68, C07K 16/18, G06F 17/30, 17/50</b>	<b>A2</b>	<b>(11) International Publication Number: WO 99/53051</b> <b>(43) International Publication Date: 21 October 1999 (21.10.99)</b>
<b>(21) International Application Number:</b> PCT/IB99/00712 <b>(22) International Filing Date:</b> 9 April 1999 (09.04.99)  <b>(30) Priority Data:</b> 09/057,719 9 April 1998 (09.04.98) US 09/069,047 28 April 1998 (28.04.98) US  <b>(71) Applicant (for all designated States except US):</b> GENSET [FR/FR]; 24, rue Royale, F-75008 Paris (FR).  <b>(72) Inventors; and</b> <b>(75) Inventors/Applicants (for US only):</b> DUMAS MILNE EDWARDS, Jean-Baptiste [FR/FR]; 8, rue Grégoire-de-Tours, F-75006 Paris (FR). DUCLERT, Aymeric [FR/FR]; 6 ter, rue Victorine, F-94100 Saint-Maur (FR). GIORDANO, Jean-Yves [FR/FR]; 12, rue Duhesme, F-75018 Paris (FR).  <b>(74) Agents:</b> MARTIN, Jean-Jacques et al.; Cabinet Regimbeau, 26, avenue Kléber, F-75116 Paris (FR).		<b>(81) Designated States:</b> AU, CA, JP, US, European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE).  <b>Published</b> <i>Without international search report and to be republished upon receipt of that report.</i>
<b>(54) Title:</b> 5' ESTS AND ENCODED HUMAN PROTEINS  <b>(57) Abstract</b>  The sequences of 5' ESTs derived from mRNAs encoding secreted proteins are disclosed. The 5' ESTs may be to obtain cDNAs and genomic DNAs corresponding to the 5' ESTs. The 5' ESTs may also be used in diagnostic, forensic, gene therapy, and chromosome mapping procedures. Upstream regulatory sequences may also be obtained using the 5' ESTs. The 5' ESTs may also be used to design expression vectors and secretion vectors.		

**FOR THE PURPOSES OF INFORMATION ONLY**

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece			TR	Turkey
BG	Bulgaria	HU	Hungary	ML	Mali	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MN	Mongolia	UA	Ukraine
BR	Brazil	IL	Israel	MR	Mauritania	UG	Uganda
BY	Belarus	IS	Iceland	MW	Malawi	US	United States of America
CA	Canada	IT	Italy	MX	Mexico	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NE	Niger	VN	Viet Nam
CG	Congo	KE	Kenya	NL	Netherlands	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NO	Norway	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	NZ	New Zealand		
CM	Cameroon			PL	Poland		
CN	China	KR	Republic of Korea	PT	Portugal		
CU	Cuba	KZ	Kazakstan	RO	Romania		
CZ	Czech Republic	LC	Saint Lucia	RU	Russian Federation		
DE	Germany	LI	Liechtenstein	SD	Sudan		
DK	Denmark	LK	Sri Lanka	SE	Sweden		
EE	Estonia	LR	Liberia	SG	Singapore		

## 5' ESTS AND ENCODED HUMAN PROTEINS

Background of the Invention

The estimated 50,000-100,000 genes scattered along the human chromosomes offer tremendous  
5 promise for the understanding, diagnosis, and treatment of human diseases. In addition, probes capable  
of specifically hybridizing to loci distributed throughout the human genome find applications in the  
construction of high resolution chromosome maps and in the identification of individuals.

In the past, the characterization of even a single human gene was a painstaking process,  
requiring years of effort. Recent developments in the areas of cloning vectors, DNA sequencing, and  
10 computer technology have merged to greatly accelerate the rate at which human genes can be isolated,  
sequenced, mapped, and characterized.

Currently, two different approaches are being pursued for identifying and characterizing the  
genes distributed along the human genome. In one approach, large fragments of genomic DNA are  
isolated, cloned, and sequenced. Potential open reading frames in these genomic sequences are  
15 identified using bioinformatics software. However, this approach entails sequencing large stretches  
of human DNA which do not encode proteins in order to find the protein encoding sequences  
scattered throughout the genome. In addition to requiring extensive sequencing, the bioinformatics  
software may mischaracterize the genomic sequences obtained, *i.e.*, labeling non-coding DNA as  
coding DNA and vice versa.

20 An alternative approach takes a more direct route to identifying and characterizing human  
genes. In this approach, complementary DNAs (cDNAs) are synthesized from isolated messenger  
RNAs (mRNAs) which encode human proteins. Using this approach, sequencing is only performed on  
DNA which is derived from protein coding portions of the genome. Often, only short stretches of the  
cDNAs are sequenced to obtain sequences called expressed sequence tags (ESTs). The ESTs may then  
25 be used to isolate or purify extended cDNAs which include sequences adjacent to the EST sequences.  
The extended cDNAs may contain all of the sequence of the EST which was used to obtain them or only  
a portion of the sequence of the EST which was used to obtain them. In addition, the extended cDNAs  
may contain the full coding sequence of the gene from which the EST was derived or, alternatively, the  
extended cDNAs may include portions of the coding sequence of the gene from which the EST was  
30 derived. It will be appreciated that there may be several extended cDNAs which include the EST  
sequence as a result of alternate splicing or the activity of alternative promoters. Alternatively, ESTs  
having partially overlapping sequences may be identified and contigs comprising the consensus  
sequences of the overlapping ESTs may be identified.

In the past, these short EST sequences were often obtained from oligo-dT primed cDNA  
35 libraries. Accordingly, they mainly corresponded to the 3' untranslated region of the mRNA. In part,  
the prevalence of EST sequences derived from the 3' end of the mRNA is a result of the fact that typical

techniques for obtaining cDNAs, are not well suited for isolating cDNA sequences derived from the 5' ends of mRNAs (Adams *et al.*, *Nature* 377:3-174, 1996, Hillier *et al.*, *Genome Res.* 6:807-828, 1996).

In addition, in those reported instances where longer cDNA sequences have been obtained, the reported sequences typically correspond to coding sequences and do not include the full 5' untranslated region (5'UTR) of the mRNA from which the cDNA is derived. Indeed, 5'UTRs have been shown to affect either the stability or translation of mRNAs. Thus, regulation of gene expression may be achieved through the use of alternative 5'UTRs as shown, for instance, for the translation of the tissue inhibitor of metalloprotease mRNA in mitogenically activated cells (Waterhouse *et al.*, *J Biol Chem.* 265:5585-9, 1990). Furthermore, modification of 5'UTR through mutation, insertion or translocation events may even be implied in pathogenesis. For instance, the fragile X syndrome, the most common cause of inherited mental retardation, is partly due to an insertion of multiple CGG trinucleotides in the 5'UTR of the fragile X mRNA resulting in the inhibition of protein synthesis via ribosome stalling (Feng *et al.*, *Science* 268:731-4, 1995). An aberrant mutation in regions of the 5'UTR known to inhibit translation of the proto-oncogene *c-myc* was shown to result in upregulation of *c-myc* protein levels in cells derived from patients with multiple myelomas (Willis *et al.*, *Curr Top Microbiol Immunol* 224:269-76, 1997). In addition, the use of oligo-dT primed cDNA libraries does not allow the isolation of complete 5'UTRs since such incomplete sequences obtained by this process may not include the first exon of the mRNA, particularly in situations where the first exon is short. Furthermore, they may not include some exons, often short ones, which are located upstream of splicing sites. Thus, there is a need to obtain sequences derived from the 5' ends of mRNAs.

While many sequences derived from human chromosomes have practical applications, approaches based on the identification and characterization of those chromosomal sequences which encode a protein product are particularly relevant to diagnostic and therapeutic uses. In some instances, the sequences used in such therapeutic or diagnostic techniques may be sequences which encode proteins which are secreted from the cell in which they are synthesized. Those sequences encoding secreted proteins as well as the secreted proteins themselves, are particularly valuable as potential therapeutic agents. Such proteins are often involved in cell to cell communication and may be responsible for producing a clinically relevant response in their target cells. In fact, several secretory proteins, including tissue plasminogen activator, G-CSF, GM-CSF, erythropoietin, human growth hormone, insulin, interferon- $\alpha$ , interferon- $\beta$ , interferon- $\gamma$ , and interleukin-2, are currently in clinical use. These proteins are used to treat a wide range of conditions, including acute myocardial infarction, acute ischemic stroke, anemia, diabetes, growth hormone deficiency, hepatitis, kidney carcinoma, chemotherapy-induced neutropenia and multiple sclerosis. For these reasons, extended cDNAs encoding secreted proteins or portions thereof represent a valuable source of therapeutic agents. Thus, there is a need for the identification and characterization of secreted proteins and the nucleic acids encoding them.

In addition to being therapeutically useful themselves, secretory proteins include short peptides, called signal peptides, at their amino termini which direct their secretion. These signal peptides are



encoded by the signal sequences located at the 5' ends of the coding sequences of genes encoding secreted proteins. These signal peptides can be used to direct the extracellular secretion of any protein to which they are operably linked. In addition, portions of the signal peptides called membrane-translocating sequences, may also be used to direct the intracellular import of a peptide or protein of interest. This may prove beneficial in gene therapy strategies in which it is desired to deliver a particular gene product to cells other than the cells in which it is produced. Signal sequences encoding signal peptides also find application in simplifying protein purification techniques. In such applications, the extracellular secretion of the desired protein greatly facilitates purification by reducing the number of undesired proteins from which the desired protein must be selected. Thus, there exists a need to identify and characterize the 5' portions of the genes for secretory proteins which encode signal peptides.

Sequences coding for non-secreted proteins may also find application as therapeutics or diagnostics. In particular, such sequences may be used to determine whether an individual is likely to express a detectable phenotype, such as a disease, as a consequence of a mutation in the coding sequence of a protein. In instances where the individual is at risk of suffering from a disease or other undesirable phenotype as a result of a mutation in such a coding sequence, the undesirable phenotype may be corrected by introducing a normal coding sequence using gene therapy. Alternatively, if the undesirable phenotype results from overexpression of the protein encoded by the coding sequence, expression of the protein may be reduced using antisense or triple helix based strategies.

The secreted or non-secreted human polypeptides encoded by the coding sequences may also be used as therapeutics by administering them directly to an individual having a condition, such as a disease, resulting from a mutation in the sequence encoding the polypeptide. In such an instance, the condition can be cured or ameliorated by administering the polypeptide to the individual.

In addition, the secreted or non-secreted human polypeptides or portions thereof may be used to generate antibodies useful in determining the tissue type or species of origin of a biological sample. The antibodies may also be used to determine the cellular localization of the secreted or non-secreted human polypeptides or the cellular localization of polypeptides which have been fused to the human polypeptides. In addition, the antibodies may also be used in immunoaffinity chromatography techniques to isolate, purify, or enrich the human polypeptide or a target polypeptide which has been fused to the human polypeptide.

Public information on the number of human genes for which the promoters and upstream regulatory regions have been identified and characterized is quite limited. In part, this may be due to the difficulty of isolating such regulatory sequences. Upstream regulatory sequences such as transcription factor binding sites are typically too short to be utilized as probes for isolating promoters from human genomic libraries. Recently, some approaches have been developed to isolate human promoters. One of them consists of making a CpG island library (Cross *et al.*, *Nature Genetics* 6: 236-244, 1994). The second consists of isolating human genomic DNA sequences containing SpeI binding sites by the use of SpeI binding protein. (Mortlock *et al.*, *Genome Res.* 6:327-335, 1996). Both of these approaches have

their limits due to a lack of specificity and of comprehensiveness. Thus, there exists a need to identify and systematically characterize the 5' portions of the genes.

The present 5' ESTs may be used to efficiently identify and isolate 5'UTRs and upstream regulatory regions which control the location, developmental stage, rate, and quantity of protein synthesis, as well as the stability of the mRNA. Once identified and characterized, these regulatory regions may be utilized in gene therapy or protein purification schemes to obtain the desired amount and locations of protein synthesis or to inhibit, reduce, or prevent the synthesis of undesirable gene products.

In addition, ESTs containing the 5' ends of protein genes may include sequences useful as probes for chromosome mapping and the identification of individuals. Thus, there is a need to identify and characterize the sequences upstream of the 5' coding sequences of genes.

#### Summary of the Invention

The present invention relates to purified, isolated, or enriched 5' ESTs which include sequences derived from the authentic 5' ends of their corresponding mRNAs. The term "corresponding mRNA" refers to the mRNA which was the template for the cDNA synthesis which produced the 5' EST. These sequences will be referred to hereinafter as "5' ESTs." The present invention also includes purified, isolated or enriched nucleic acids comprising contigs assembled by determining a consensus sequences from a plurality of ESTs containing overlapping sequences. These contigs will be referred to herein as "consensus contigated 5'ESTs."

As used herein, the term "purified" does not require absolute purity; rather, it is intended as a relative definition. Individual 5' EST clones isolated from a cDNA library have been conventionally purified to electrophoretic homogeneity. The sequences obtained from these clones could not be obtained directly either from the library or from total human DNA. The cDNA clones are not naturally occurring as such, but rather are obtained via manipulation of a partially purified naturally occurring substance (messenger RNA). The conversion of mRNA into a cDNA library involves the creation of a synthetic substance (cDNA) and pure individual cDNA clones can be isolated from the synthetic library by clonal selection. Thus, creating a cDNA library from messenger RNA and subsequently isolating individual clones from that library results in an approximately  $10^4$ - $10^6$  fold purification of the native message. Purification of starting material or natural material to at least one order of magnitude, preferably two or three orders, and more preferably four or five orders of magnitude is expressly contemplated.

As used herein, the term "isolated" requires that the material be removed from its original environment (e.g., the natural environment if it is naturally occurring). For example, a naturally-occurring polynucleotide present in a living animal is not isolated, but the same polynucleotide, separated from some or all of the coexisting materials in the natural system, is isolated.

As used herein, the term "recombinant" means that the 5' EST is adjacent to "backbone" nucleic acid to which it is not adjacent in its natural environment. Additionally, to be "enriched" the 5' ESTs will

represent 5% or more of the number of nucleic acid inserts in a population of nucleic acid backbone molecules. Backbone molecules according to the present invention include nucleic acids such as expression vectors, self-replicating nucleic acids, viruses, integrating nucleic acids, and other vectors or nucleic acids used to maintain or manipulate a nucleic acid insert of interest. Preferably, the enriched 5'

- 5 ESTs represent 15% or more of the number of nucleic acid inserts in the population of recombinant backbone molecules. More preferably, the enriched 5' ESTs represent 50% or more of the number of nucleic acid inserts in the population of recombinant backbone molecules. In a highly preferred embodiment, the enriched 5' ESTs represent 90% or more of the number of nucleic acid inserts in the population of recombinant backbone molecules.

- 10 "Stringent," "moderate," and "low" hybridization conditions are as defined below.

The term "polypeptide" refers to a polymer of amino acids without regard to the length of the polymer; thus, peptides, oligopeptides, and proteins are included within the definition of polypeptide. This term also does not specify or exclude post-expression modifications of polypeptides, for example, polypeptides which include the covalent attachment of glycosyl groups, acetyl groups,

- 15 phosphate groups, lipid groups and the like are expressly encompassed by the term polypeptide. Also included within the definition are polypeptides which contain one or more analogs of an amino acid (including, for example, non-naturally occurring amino acids, amino acids which only occur naturally in an unrelated biological system, modified amino acids from mammalian systems etc.), polypeptides with substituted linkages, as well as other modifications known in the art, both naturally  
20 occurring and non-naturally occurring.

- As used interchangeably herein, the terms "nucleic acids," "oligonucleotides," and "polynucleotides" include RNA, DNA, or RNA/DNA hybrid sequences of more than one nucleotide in either single chain or duplex form. The term "nucleotide" as used herein as an adjective to describe molecules comprising RNA, DNA, or RNA/DNA hybrid sequences of any length in single-  
25 stranded or duplex form. The term "nucleotide" is also used herein as a noun to refer to individual nucleotides or varieties of nucleotides, meaning a molecule, or individual unit in a larger nucleic acid molecule, comprising a purine or pyrimidine, a ribose or deoxyribose sugar moiety, and a phosphate group, or phosphodiester linkage in the case of nucleotides within an oligonucleotide or polynucleotide. Although the term "nucleotide" is also used herein to encompass "modified  
30 nucleotides" which comprise at least one modifications (a) an alternative linking group, (b) an analogous form of purine, (c) an analogous form of pyrimidine, or (d) an analogous sugar, for examples of analogous linking groups, purine, pyrimidines, and sugars see for example PCT publication No. WO 95/04064. The polynucleotide sequences of the invention may be prepared by any known method, including synthetic, recombinant, *ex vivo* generation, or a combination thereof,  
35 as well as utilizing any purification methods known in the art.

The terms "base paired" and "Watson & Crick base paired" are used interchangeably herein to refer to nucleotides which can be hydrogen bonded to one another by virtue of their sequence

identities in a manner like that found in double-helical DNA with thymine or uracil residues linked to adenine residues by two hydrogen bonds and cytosine and guanine residues linked by three hydrogen bonds (See Stryer, L., *Biochemistry*, 4<sup>th</sup> edition, 1995).

The terms "complementary" or "complement thereof" are used herein to refer to the  
5 sequences of polynucleotides which are capable of forming Watson & Crick base pairing with another specified polynucleotide throughout the entirety of the complementary region. For the purpose of the present invention, a first polynucleotide is deemed to be complementary to a second polynucleotide when each base in the first polynucleotide is paired with its complementary base. Complementary bases are, generally, A and T (or A and U), or C and G. "Complement" is used  
10 herein as a synonym from "complementary polynucleotide," "complementary nucleic acid" and "complementary nucleotide sequence". These terms are applied to pairs of polynucleotides based solely upon their sequences and not any particular set of conditions under which the two polynucleotides would actually bind. Preferably, a "complementary" sequence is a sequence which an A at each position where there is a T on the opposite strand, a T at each position where there is an A on  
15 the opposite strand, a G at each position where there is a C on the opposite strand and a C at each position where there is a G on the opposite strand.

Thus, 5' ESTs in cDNA libraries in which one or more 5' ESTs make up 5% or more of the number of nucleic acid inserts in the backbone molecules are "enriched recombinant 5' ESTs" as defined herein. Likewise, 5' ESTs in a population of plasmids in which one or more 5' ESTs of the present  
20 invention have been inserted such that they represent 5% or more of the number of inserts in the plasmid backbone are "enriched recombinant 5' ESTs" as defined herein. However, 5' ESTs in cDNA libraries in which 5' ESTs constitute less than 5% of the number of nucleic acid inserts in the population of backbone molecules, such as libraries in which backbone molecules having a 5' EST insert are extremely rare, are not "enriched recombinant 5' ESTs."

25 In some embodiments, the present invention relates to 5' ESTs which are derived from genes encoding secreted proteins. As used herein, a "secreted" protein is one which, when expressed in a suitable host cell, is transported across or through a membrane, including transport as a result of signal peptides in its amino acid sequence. "Secreted" proteins include without limitation proteins secreted wholly (e.g. soluble proteins), or partially (e.g. receptors) from the cell in which they are expressed.  
30 "Secreted" proteins also include without limitation proteins which are transported across the membrane of the endoplasmic reticulum.

Such 5' ESTs include nucleic acid sequences, called signal sequences, which encode signal peptides which direct the extracellular secretion of the proteins encoded by the genes from which the 5' ESTs are derived. Generally, the signal peptides are located at the amino termini of secreted proteins.

35 Secreted proteins are translated by ribosomes associated with the "rough" endoplasmic reticulum. Generally, secreted proteins are co-translationally transferred to the membrane of the endoplasmic reticulum. Association of the ribosome with the endoplasmic reticulum during translation

of secreted proteins is mediated by the signal peptide. The signal peptide is typically cleaved following its co-translational entry into the endoplasmic reticulum. After delivery to the endoplasmic reticulum, secreted proteins may proceed through the Golgi apparatus. In the Golgi apparatus, the proteins may undergo post-translational modification before entering secretory vesicles which transport them across  
5 the cell membrane.

The 5' ESTs of the present invention have several important applications. For example, they may be used to obtain and express cDNA clones which include the full protein coding sequences of the corresponding gene products, including the authentic translation start sites derived from the 5' ends of the coding sequences of the mRNAs from which the 5' ESTs are derived. These cDNAs will be referred  
10 to hereinafter as "full-length cDNAs." These cDNAs may also include DNA derived from mRNA sequences upstream of the translation start site. The full-length cDNA sequences may be used to express the proteins corresponding to the 5' ESTs. As discussed above, secreted proteins and non-secreted proteins may be therapeutically important. Thus, the proteins expressed from the cDNAs may be useful in treating and controlling a variety of human conditions. The 5' ESTs may also be used to obtain the  
15 corresponding genomic DNA. The term "corresponding genomic DNA" refers to the genomic DNA which encodes the mRNA from which the 5' EST was derived.

Alternatively, the 5' ESTs may be used to obtain and express extended cDNAs encoding portions of the protein. In the case of secreted proteins, the portions may comprise the signal peptides of the secreted proteins or the mature proteins generated when the signal peptide is cleaved off.

20 The present invention includes isolated, purified, or enriched "EST-related nucleic acids." The terms "isolated," "purified" or "enriched" have the meanings provided above. As used herein, the term "EST-related nucleic acids" means the nucleic acids of SEQ ID NOs. 24-811 and 1600-1622, extended cDNAs obtainable using the nucleic acids of SEQ ID NOs. 24-811 and 1600-1622, full-length cDNAs obtainable using the nucleic acids of SEQ ID NOs. 24-811 and 1600-1622 or genomic DNAs obtainable  
25 using the nucleic acids of SEQ ID NOs. 24-811 and 1600-1622. The present invention also includes the sequences complementary to the EST-related nucleic acids.

The present invention also includes isolated, purified, or enriched "fragments of EST-related nucleic acids." The terms "isolated," "purified" and "enriched" have the meanings described above. As used herein the term "fragments of EST-related nucleic acids" means fragments comprising at least 10,  
30 12, 15, 18, 20, 23, 25, 28, 30, 35, 40, 50, 75, 100, 200, 300, 500, or 1000 consecutive nucleotides of the EST-related nucleic acids to the extent that fragments of these lengths are consistent with the lengths of the particular EST-related nucleic acids being referenced. In particular, fragments of EST-related nucleic acids refer to "polynucleotides described in Table II," "polynucleotides described in Table III," and "polynucleotides described in Table IV." The present invention also includes the sequences  
35 complementary to the fragments of the EST-related nucleic acids.

The present invention also includes isolated, purified, or enriched "positional segments of EST-related nucleic acids." As used herein, the term "positional segments of EST-related nucleic acids"

includes segments comprising nucleotides 1-25, 26-50, 51-75, 76-100, 101-125, 126-150, 151-175, 176-200, 201-225, 226-250, 251-300, 301-325, 326-350, 351-375, 376-400, 401-425, 426-450, 451-475, 476-500, 501-525, 526-550, 551-575, 576-600 and 601-the terminal nucleotide of the EST-related nucleic acids to the extent that such nucleotide positions are consistent with the lengths of the particular

5 EST-related nucleic acids being referenced. The term "positional segments of EST-related nucleic acids" also includes segments comprising nucleotides 1-50, 51-100, 101-150, 151-200, 201-250, 251-300, 301-350, 351-400, 401-450, 450-500, 501-550, 551-600 or 601-the terminal nucleotide of the EST-related nucleic acids to the extent that such nucleotide positions are consistent with the lengths of the particular

10 EST-related nucleic acids being referenced. The term "positional segments of EST-related nucleic acids" also includes segments comprising nucleotides 1-100, 101-200, 201-300, 301-400, 501-500, 500-600, or 601-the terminal nucleotide of the EST-related nucleic acids to the extent that such nucleotide positions are consistent with the lengths of the particular EST-related nucleic acids being referenced. In addition, the term "positional segments of EST-related nucleic acids" includes segments comprising nucleotides 1-200, 201-400, 400-600, or 601-the terminal nucleotide of the EST-related nucleic acids to

15 the extent that such nucleotide positions are consistent with the lengths of the particular EST-related nucleic acids being referenced. The present invention also includes the sequences complementary to the positional segments of EST-related nucleic acids.

The present invention also includes isolated, purified, or enriched "fragments of positional segments of EST-related nucleic acids." As used herein, the term "fragments of positional segments of

20 EST-related nucleic acids" refers to fragments comprising at least 10, 15, 18, 20, 23, 25, 28, 30, 35, 40, 50, 75, 100, 150, or 200 consecutive nucleotides of the positional segments of EST-related nucleic acids. The present invention also includes the sequences complementary to the fragments of positional segments of EST-related nucleic acids.

The present invention also includes isolated or purified "EST-related polypeptides." As used

25 herein, the term "EST-related polypeptides" means the polypeptides encoded by the EST-related nucleic acids, including the polypeptides of SEQ ID NOs. 812-1599.

The present invention also includes isolated or purified "fragments of EST-related polypeptides." As used herein, the term "fragments of EST-related polypeptides" means fragments comprising at least 5, 10, 15, 20, 25, 30, 35, 40, 50, 75, 100, or 150 consecutive amino acids of an EST-

30 related polypeptide to the extent that fragments of these lengths are consistent with the lengths of the particular EST-related polypeptides being referenced. In particular, fragments of EST-related polypeptides refer to polypeptides encoded by "polynucleotides described in Table II," "polynucleotides described in Table III," and "polynucleotides described in Table IV."

The present invention also includes isolated or purified "positional segments of EST-related

35 polypeptides." As used herein, the term "positional segments of EST-related polypeptides" includes polypeptides comprising amino acid residues 1-25, 26-50, 51-75, 76-100, 101-125, 126-150, 151-175, 176-200, or 201-the C-terminal amino acid of the EST-related polypeptides to the extent that such amino

acid residues are consistent with the lengths of the particular EST-related polypeptides being referenced. The term "positional segments of EST-related polypeptides" also includes segments comprising amino acid residues 1-50, 51-100, 101-150, 151-200 or 201-the C-terminal amino acid of the EST-related polypeptides to the extent that such amino acid residues are consistent with the lengths of the particular EST-related polypeptides being referenced. The term "positional segments of EST-related polypeptides" also includes segments comprising amino acids 1-100 or 101-200 of the EST-related polypeptides to the extent that such amino acid residues are consistent with the lengths of particular EST-related polypeptides being referenced. In addition, the term "positional segments of EST-related polypeptides" includes segments comprising amino acid residues 1-200 or 201-the C-terminal amino acid of the EST-related polypeptides to the extent that amino acid residues are consistent with the lengths of the particular EST-related polypeptides being referenced.

The present invention also includes isolated or purified "fragments of positional segments of EST-related polypeptides." As used herein, the term "fragments of positional segments of EST-related polypeptides" means fragments comprising at least 5, 10, 15, 20, 25, 30, 35, 40, 50, 75, 100, or 150 consecutive amino acids of positional segments of EST-related polypeptides to the extent that fragments of these lengths are consistent with the lengths of the particular EST-related polypeptides being referenced.

The present invention also includes antibodies which specifically recognize the EST-related polypeptides, fragments of EST-related polypeptides, positional segments of EST-related polypeptides, or fragments of positional segments of EST-related polypeptides. In the case of secreted proteins, such as those of SEQ ID NOs. 1554-1580 antibodies which specifically recognize the mature protein generated when the signal peptide is cleaved may also be obtained as described below. Similarly, antibodies which specifically recognize the signal peptides of SEQ ID NOs. 812-1516 or 1554-1580 may also be obtained.

In some embodiments and in the case of secreted proteins, the EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of EST-related nucleic acids, or fragments of positional segments of nucleic acids include a signal sequence. In other embodiments, the EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of EST-related nucleic acids, or fragments of positional segments of nucleic acids may include the full coding sequence for the protein or, in the case of secreted proteins, the full coding sequence of the mature protein (*i.e.* the protein generated when the signal polypeptide is cleaved off). In addition, the EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of EST-related nucleic acids, or fragments of positional segments of nucleic acids may include regulatory regions upstream of the translation start site or downstream of the stop codon which control the amount, location, or developmental stage of gene expression.

As discussed above, both secreted and non-secreted human proteins may be therapeutically important. Thus, the proteins expressed from the EST-related nucleic acids, fragments of EST-related

nucleic acids, positional segments of EST-related nucleic acids, or fragments of positional segments of nucleic acids may be useful in treating or controlling a variety of human conditions.

The EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of EST-related nucleic acids, or fragments of positional segments of nucleic acids may be used in forensic  
5 procedures to identify individuals or in diagnostic procedures to identify individuals having genetic diseases resulting from abnormal gene expression. In addition, the EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of EST-related nucleic acids, or fragments of positional segments of nucleic acids are useful for constructing a high resolution map of the human chromosomes.

10 The present invention also relates to secretion vectors capable of directing the secretion of a protein of interest. Such vectors may be used in gene therapy strategies in which it is desired to produce a gene product in one cell which is to be delivered to another location in the body. Secretion vectors may also facilitate the purification of desired proteins.

The present invention also relates to expression vectors capable of directing the expression of an  
15 inserted gene in a desired spatial or temporal manner or at a desired level. Such vectors may include sequences upstream of the EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of EST-related nucleic acids, or fragments of positional segments of nucleic acids, such as promoters or upstream regulatory sequences.

The present invention also comprises fusion vectors for making chimeric polypeptides  
20 comprising a first polypeptide and a second polypeptide. Such vectors are useful for determining the cellular localization of the chimeric polypeptides or for isolating, purifying or enriching the chimeric polypeptides.

The EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of EST-related nucleic acids, or fragments of positional segments of nucleic acids may also be used for  
25 gene therapy to control or treat genetic diseases. In the case of secreted proteins, signal peptides may be fused to heterologous proteins to direct their extracellular secretion.

Bacterial clones containing Bluescript plasmids having inserts containing the sequence of the non-aligned 5'ESTs, also referred to as singletons, and sequences of the 5'ESTs which were aligned to yield consensus contigated 5' ESTs are presently stored at 80°C in 4% (v/v) glycerol in the inventor's  
30 laboratories under internal designations. The non-aligned 5'ESTs are those which comprise a single EST from a single tissue in the listing of Table V. The inserts may be recovered from the stored materials by growing the appropriate clones on a suitable medium. The Bluescript DNA can then be isolated using plasmid isolation procedures familiar to those skilled in the art such as alkaline lysis minipreps or large scale alkaline lysis plasmid isolation procedures. If desired the plasmid DNA may be  
35 further enriched by centrifugation on a cesium chloride gradient, size exclusion chromatography, or anion exchange chromatography. The plasmid DNA obtained using these procedures may then be manipulated using standard cloning techniques familiar to those skilled in the art. Alternatively, a PCR



can be performed with primers designed at both ends of the inserted EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of EST-related nucleic acids, or fragments of positional segments of nucleic acids. The PCR product which corresponds to the EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of EST-related nucleic acids, or  
5 fragments of positional segments of nucleic acids can then be manipulated using standard cloning techniques familiar to those skilled in the art.

One embodiment of the present invention is a purified nucleic acid comprising a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 and sequences complementary to the sequences of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622.

10 Another embodiment of the present invention is a purified nucleic acid comprising at least 10, 12, 15, 18, 20, 23, 25, 28, 30, 35, 40, 50, 75, 100, 200, 300, 500, or 1000 consecutive nucleotides, to the extent that fragments of these lengths are consistent with the specific sequence, of a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 and sequences complementary to the sequences of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622.

15 A further embodiment of the present invention is a purified nucleic acid comprising the coding sequence of a sequence selected from the group consisting of SEQ ID NOs. 24-811.

Yet another embodiment of the present invention is a purified nucleic acid comprising the full coding sequences of a sequence selected from the group consisting of SEQ ID NOs. 766-792 wherein the full coding sequence comprises the sequence encoding the signal peptide and the  
20 sequence encoding the mature protein.

Still another embodiment of the present invention is a purified nucleic acid comprising a contiguous span of a sequence selected from the group consisting of SEQ ID NOs. 766-792 which encodes the mature protein.

Another embodiment of the present invention is a purified nucleic acid comprising a  
25 contiguous span of a sequence selected from the group consisting of SEQ ID NOs. 24-728 and 766-792 which encodes the signal peptide.

Another embodiment of the present invention is a purified nucleic acid encoding a polypeptide comprising a sequence selected from the group consisting of the sequences of SEQ ID NOs. 812-1599.

30 Another embodiment of the present invention is a purified nucleic acid encoding a polypeptide comprising a sequence selected from the group consisting of the sequences of SEQ ID NOs. 1554-1580.

Another embodiment of the present invention is a purified nucleic acid encoding a polypeptide comprising a mature protein included in a sequence selected from the group consisting of  
35 the sequences of SEQ ID NOs. 1554-1580.

Another embodiment of the present invention is a purified nucleic acid encoding a polypeptide comprising a signal peptide included in a sequence selected from the group consisting of the sequences of SEQ ID NOs. 812-1516 and 1554-1580.

Another embodiment of the present invention is a purified nucleic acid at least 30, 35, 40, 50, 75, 100, 200, 300, 500 or 1000 nucleotides in length which hybridizes under stringent conditions to a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 and sequences complementary to the sequences of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622.

Another embodiment of the present invention is a purified or isolated polypeptide comprising a sequence selected from the group consisting of the sequences of SEQ ID NOs. 812-1599.

Another embodiment of the present invention is a purified or isolated polypeptide comprising a sequence selected from the group consisting of SEQ ID NOs. 1554-1580.

Another embodiment of the present invention is a purified or isolated polypeptide comprising a mature protein of a polypeptide selected from the group consisting of SEQ ID NOs. 1554-1580.

Another embodiment of the present invention is a purified or isolated polypeptide comprising a signal peptide of a sequence selected from the group consisting of the polypeptides of SEQ ID NOs. 812-1516 and 1554-1580.

Another embodiment of the present invention is a purified or isolated polypeptide comprising at least 12, 15, 18, 20, 23, 25, 28, 30, 35, 40, 50, 75, 100, 200, 300, 500, or 1000 consecutive amino acids, to the extent that fragments of these lengths are consistent with the specific sequence, of a sequence selected from the group consisting of the sequences of SEQ ID NOs. 812-1599.

Another embodiment of the present invention is a method of making a cDNA comprising the steps of contacting a collection of mRNA molecules from human cells with a primer comprising at least 12, 15, 18, 20, 23, 25, 28, 30, 35, 40, or 50 consecutive nucleotides of a sequence selected from the group consisting of the sequences complementary to SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622, hybridizing said primer to an mRNA in said collection that encodes said protein reverse transcribing said hybridized primer to make a first cDNA strand from said mRNA, making a second cDNA strand complementary to said first cDNA strand and isolating the resulting cDNA encoding said protein comprising said first cDNA strand and said second cDNA strand.

Another embodiment of the present invention is a purified cDNA obtainable by the method of the preceding paragraph.

In one aspect of this embodiment, the cDNA encodes at least a portion of a human polypeptide.

Another embodiment of the present invention is a method of making a cDNA comprising the steps of obtaining a cDNA comprising a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622, contacting said cDNA with a detectable probe comprising at least 12, 15, 18, 20, 23, 25, 28, 30, 35, 40, or 50 consecutive nucleotides of a sequence

selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 and the sequences complementary to SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 under conditions which permit said probe to hybridize to said cDNA, identifying a cDNA which hybridizes to said detectable probe, and isolating said cDNA which hybridizes to said probe.

5 Another embodiment of the present invention is a purified cDNA obtainable by the method of the preceding paragraph.

In one aspect of this embodiment, the cDNA encodes at least a portion of a human polypeptide.

Another embodiment of the present invention is a method of making a cDNA comprising the  
10 steps of contacting a collection of mRNA molecules from human cells with a first primer capable of hybridizing to the polyA tail of said mRNA, hybridizing said first primer to said polyA tail, reverse transcribing said mRNA to make a first cDNA strand, making a second cDNA strand complementary to said first cDNA strand using at least one primer comprising at least 12, 15, 18, 20, 23, 25, 28, 30, 35, 40, or 50 consecutive nucleotides of a sequence selected from the group consisting of SEQ ID  
15 NOs. 24-811 and SEQ ID NOs. 1600-1622, and isolating the resulting cDNA comprising said first cDNA strand and said second cDNA strand.

Another embodiment of the present invention is a purified cDNA obtainable by the method of the preceding paragraph.

In one aspect of this embodiment, said cDNA encodes at least a portion of a human  
20 polypeptide.

In another aspect of the preceding method the second cDNA strand is made by contacting said first cDNA strand with a first pair of primers, said first pair of primers comprising a second primer comprising at least 12, 15, 18, 20, 23, 25, 28, 30, 35, 40, or 50 consecutive nucleotides of a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622  
25 and a third primer having a sequence therein which is included within the sequence of said first primer, performing a first polymerase chain reaction with said first pair of primers to generate a first PCR product, contacting said first PCR product with a second pair of primers, said second pair of primers comprising a fourth primer, said fourth primer comprising at least 12, 15, 18, 20, 23, 25, 28, 30, 35, 40, or 50 consecutive nucleotides of said sequence selected from the group consisting of SEQ  
30 ID NOs. 24-811 and SEQ ID NOs. 1600-1622, and a fifth primer, wherein said fourth and fifth hybridize to sequences within said first PCR product, and performing a second polymerase chain reaction, thereby generating a second PCR product.

One aspect of this embodiment is a purified cDNA obtainable by the method of the preceding paragraph.

35 In another aspect of this embodiment, said cDNA encodes at least a portion of a human polypeptide.

Alternatively, the second cDNA strand may be made by contacting said first cDNA strand with a second primer comprising at least 12, 15, 18, 20, 23, 25, 28, 30, 35, 40, or 50 consecutive nucleotides of a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622, hybridizing said second primer to said first strand cDNA, and extending said  
5 hybridized second primer to generate said second cDNA strand.

One aspect of the above embodiment is a purified cDNA obtainable by the method of the preceding paragraph.

In a further aspect of this embodiment said cDNA encodes at least a portion of a human polypeptide.

10 Another embodiment of the present invention is a method of making a polypeptide comprising the steps of obtaining a cDNA which encodes a polypeptide encoded by a nucleic acid comprising a sequence selected from the group consisting of SEQ ID NOs. 24-811 or a cDNA which encodes a polypeptide comprising at least 6, 8, 10, 12, 15, 18, 20, 23, 25, 28, 30, 35, 40, or 50 consecutive amino acids of a polypeptide encoded by a sequence selected from the group consisting  
15 of SEQ ID NOs. 24-811, inserting said cDNA in an expression vector such that said cDNA is operably linked to a promoter, introducing said expression vector into a host cell whereby said host cell produces the protein encoded by said cDNA, and isolating said protein.

Another aspect of this embodiment is an isolated protein obtainable by the method of the preceding paragraph.

20 Another embodiment of the present invention is a method of obtaining a promoter DNA comprising the steps of obtaining genomic DNA located upstream of a nucleic acid comprising a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 and the sequences complementary to the sequences of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622, screening said genomic DNA to identify a promoter capable of directing transcription  
25 initiation, and isolating said DNA comprising said identified promoter.

In one aspect of this embodiment, said obtaining step comprises walking from genomic DNA comprising a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 and the sequences complementary to SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622. In another aspect of this embodiment, said screening step comprises inserting genomic DNA located  
30 upstream of a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 and the sequences complementary to SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 into a promoter reporter vector. For example, said screening step may comprise identifying motifs in genomic DNA located upstream of a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 and the sequences complementary to SEQ ID NOs.  
35 24-811 and SEQ ID NOs. 1600-1622 which are transcription factor binding sites or transcription start sites.

Another embodiment of the present invention is a isolated promoter obtainable by the method of the paragraph above.

Another embodiment of the present invention is the inclusion of at least one sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622, the  
5 sequences complementary to the sequences of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 and fragments comprising at least 12, 15, 18, 20, 23, 25, 28, 30, 35, 40, 50, or 100 consecutive nucleotides of said sequence in an array of discrete ESTs or fragments thereof of at least 12, 15, 18, 20, 23, 25, 28, 30, 35, 40, 50, or 100 nucleotides in length. In some aspects of this embodiment, the array includes at least two sequences selected from the group consisting of SEQ ID NOs. 24-811 and  
10 SEQ ID NOs. 1600-1622, the sequences complementary to the sequences of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622, and fragments comprising at least 12, 15, 18, 20, 23, 25, 28, 30, 35, 40, 50, or 100 consecutive nucleotides of said sequences. In another aspect of this embodiment, the array includes at least one, three, five, ten, fifteen, or twenty sequences selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622, the sequences complementary to  
15 the sequences of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 and fragments comprising at least 12, 15, 18, 20, 23, 25, 28, 30, 35, 40, 50, or 100 consecutive nucleotides of said sequences.

Another embodiment of the present invention is an enriched population of recombinant nucleic acids, said recombinant nucleic acids comprising an insert nucleic acid and a backbone nucleic acid, wherein at least 0.01%, 0.05%, 0.1%, 0.5%, 1%, 2%, 5%, 10%, or 20% of said insert  
20 nucleic acids in said population comprise a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 and the sequences complementary to SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622.

Another embodiment of the present invention is a purified or isolated antibody capable of specifically binding to a polypeptide comprising a sequence selected from the group consisting of  
25 SEQ ID NOs. 812-1599.

Another embodiment of the present invention is a purified or isolated antibody capable of specifically binding to a polypeptide comprising at least 6, 8, 10, 12, 15, 18, 20, 23, 25, 28, 30, 35, 40, or 50 consecutive amino acids of a sequence selected from the group consisting of SEQ ID NOs. 812-1599.

30 Yet, another embodiment of the present invention is an antibody composition capable of selectively binding to an epitope-containing fragment of a polypeptide comprising a contiguous span of at least 8, 10, 12, 15, 18, 20, 23, 25, 28, 30, 35, 40, or 50 amino acids of any of SEQ ID NOs. 812-1599, wherein said antibody is polyclonal or monoclonal.

Another embodiment of the present invention is a computer readable medium having stored  
35 thereon a sequence selected from the group consisting of a nucleic acid code of SEQ ID NOs. 24-811 and 1600-1622 and a polypeptide code of SEQ ID NOs. 812-1599.

Another embodiment of the present invention is a computer system comprising a processor and a data storage device wherein said data storage device has stored thereon a sequence selected from the group consisting of a nucleic acid code of SEQID NOs. 24-811 and 1600-1622 and a polypeptide code of SEQ ID NOs. 812-1599. In one aspect of this embodiment the computer system  
5 further comprises a sequence comparer and a data storage device having reference sequences stored thereon. For example, the sequence comparer may comprise a computer program which indicates polymorphisms. In another aspect of this embodiment, the computer system further comprises an identifier which identifies features in said sequence.

Another embodiment of the present invention is a method for comparing a first sequence to a  
10 reference sequence wherein said first sequence is selected from the group consisting of a nucleic acid code of SEQID NOs. 24-811 and 1600-1622 and a polypeptide code of SEQ ID NOs. 812-1599 comprising the steps of reading said first sequence and said reference sequence through use of a computer program which compares sequences and determining differences between said first sequence and said reference sequence with said computer program. In some aspects of this embodiment, said step  
15 of determining differences between the first sequence and the reference sequence comprises identifying polymorphisms.

Another embodiment of the present invention is a method for identifying a feature in a sequence selected from the group consisting of a nucleic acid code of SEQID NOs. 24-811 and 1600-1622 and a polypeptide code of SEQ ID NOs. 812-1599 comprising the steps of reading said  
20 sequence through the use of a computer program which identifies features in sequences and identifying features in said sequence with said computer program.

Another embodiment of the present invention is a vector comprising a nucleic acid according to any one of the nucleic acids described above.

Another embodiment of the present invention is a host cell containing the above vector.

25 Another embodiment of the present invention is a method of making any of the nucleic acids described above comprising the steps of introducing said nucleic acid into a host cell such that said nucleic acid is present in multiple copies in each host cell and isolating said nucleic acid from said host cell.

Another embodiment of the present invention is a method of making a nucleic acid of any of  
30 the nucleic acids described above comprising the step of sequentially linking together the nucleotides in said nucleic acids.

Another embodiment of the present invention is a method of making any of the polypeptides described above wherein said polypeptides is 150 amino acids in length or less comprising the step of sequentially linking together the amino acids in said polypeptide.

35 Another embodiment of the present invention is a method of making any of the polypeptides described above wherein said polypeptides is 120 amino acids in length or less comprising the step of sequentially linking together the amino acids in said polypeptides.

### Brief Description of the Drawings

Figure 1 is a summary of a procedure for obtaining cDNAs which have been selected to include the 5' ends of the mRNAs from which they derived. In the first step (1), the cap of intact mRNAs is oxidized to be chemically ligated to an oligonucleotide tag. In the second step (2), a reverse transcription is performed using random primers to generate a first cDNA strand. In the third step (3), mRNAs are eliminated and the second strand synthesis is carried out using a primer contained in the oligonucleotide tag.

Figure 2 is an analysis of the 43 amino terminal amino acids of all human SwissProt proteins to determine the frequency of false positives and false negatives using the techniques for signal peptide identification described herein.

Figure 3 summarizes a general method used to clone and sequence extended cDNAs containing sequences adjacent to 5'ESTs.

Figure 4 provides a schematic description of the promoters isolated and the way they are assembled with the corresponding 5' tags.

Figure 5 describes the transcription factor binding sites present in each of the promoters of Figure 4.

Figure 6 is a block diagram of an exemplary computer system.

Figure 7 is a flow diagram illustrating one embodiment of a process for comparing a new nucleotide or protein sequence with a database of sequences in order to determine the homology levels between the new sequence and the sequences in the database.

Figure 8 is a flow diagram illustrating one embodiment of a process in a computer for determining whether two sequences are homologous.

Figure 9 is a flow diagram illustrating one embodiment of an identifier process for detecting the presence of a feature in a sequence.

Figure 10 is a table with all of the parameters that can be used for each step of extended cDNA analysis.

### Detailed Description of the Preferred Embodiment

#### **I. Obtaining 5'ESTs from cDNA libraries including the 5'Ends of their Corresponding mRNAs**

The 5' ESTs of the present invention were obtained from cDNA libraries including cDNAs which include the 5'end of their corresponding mRNAs. The general method used to obtain such cDNA libraries is described in Examples 1 to 5.

#### **EXAMPLE 1**

##### Preparation of mRNA

Total human RNAs or polyA<sup>+</sup> RNAs derived from 29 different tissues were respectively purchased from LABIMO and CLONTECH and used to generate 44 cDNA libraries as described below.

The purchased RNA had been isolated from cells or tissues using acid guanidium thiocyanate-phenol-chloroform extraction (Chomczynski and Sacchi, *Analytical Biochemistry* 162:156-159, 1987). PolyA<sup>+</sup> RNA was isolated from total RNA (LABIMO) by two passes of oligo dT chromatography, as described by Aviv and Leder, *Proc. Natl. Acad. Sci. USA* 69:1408-1412, 1972) in order to eliminate ribosomal  
5 RNA.

The quality and the integrity of the polyA<sup>+</sup> RNAs were checked. Northern blots hybridized with a globin probe were used to confirm that the mRNAs were not degraded. Contamination of the polyA<sup>+</sup> mRNAs by ribosomal sequences was checked using Northern blots and a probe derived from the sequence of the 28S rRNA. Preparations of mRNAs with less than 5% of rRNAs were used in  
10 library construction. To avoid constructing libraries with RNAs contaminated by exogenous sequences (prokaryotic or fungal), the presence of bacterial 16S ribosomal sequences or of two highly expressed fungal mRNAs was examined using PCR.

## EXAMPLE 2

### 15 Methods for Obtaining mRNAs having Intact 5' Ends

Following preparation of the mRNAs from various tissues as described above, selection of mRNA with intact 5' ends and specific attachment of an oligonucleotide tag to the 5' end of such mRNA was performed using either a chemical or enzymatic approach. Both techniques takes advantage of the presence of the "cap" structure, which characterizes the 5' end of intact mRNAs and which comprises a  
20 guanosine generally methylated once, at the 7 position. The chemical approach is illustrated in Figure 1.

The chemical modification approach involves the optional elimination of the 2', 3'-cis diol of the 3' terminal ribose, the oxidation of the 2', 3', -cis diol of the ribose linked to the cap of the 5' ends of the mRNAs into a dialdehyde, and the coupling of the such obtained dialdehyde to a derivatized oligonucleotide tag. Further detail regarding the chemical approaches for obtaining mRNAs having  
25 intact 5' ends are disclosed in International Application No. WO96/34981, published November 7, 1996.

The enzymatic approach for ligating the oligonucleotide tag to the 5' ends of mRNAs with intact 5' ends involves the removal of the phosphate groups present on the 5' ends of uncapped incomplete mRNAs, the subsequent decapping of mRNAs with intact 5' ends and the ligation of the phosphate present at the 5' end of the decapped mRNA to an oligonucleotide tag. Further detail regarding the  
30 enzymatic approaches for obtaining mRNAs having intact 5' ends are disclosed in Dumas Milne Edwards J.B. (Doctoral Thesis of Paris VI University, Le clonage des ADNc complets: difficultes et perspectives nouvelles. Apports pour l'etude de la regulation de l'expression de la tryptophane hydroxylase de rat, 20 Dec. 1993), EP0 625572 and Kato *et al.*, *Gene* 150:243-250 (1994).

In either the chemical or the enzymatic approach, the oligonucleotide tag has a restriction  
35 enzyme site (e.g. EcoRI sites) therein to facilitate later cloning procedures. Following attachment of the oligonucleotide tag to the mRNA, the integrity of the mRNA was then examined by performing a Northern blot using a probe complementary to the oligonucleotide tag.



**EXAMPLE 3****cDNA Synthesis Using mRNA Templates Having Intact 5' Ends**

For the mRNAs joined to oligonucleotide tags, first strand cDNA synthesis was performed using  
5 a reverse transcriptase with random nonamers as primers. In order to protect internal EcoRI sites in the cDNA from digestion at later steps in the procedure, methylated dCTP was used for first strand synthesis. After removal of mRNA by an alkaline hydrolysis, the first strand of cDNA was precipitated using isopropanol in order to eliminate residual primers.

The second strand of the cDNA was synthesized with a Klenow fragment using a primer  
10 corresponding to the 5' end of the ligated oligonucleotide. Methylated dCTP was also used for second strand synthesis in order to protect internal EcoRI sites in the cDNA from digestion during the cloning process.

**EXAMPLE 4****Cloning of cDNAs derived from mRNA with intact 5' ends into BlueScript**

Following second strand synthesis, the ends of the cDNA were blunted with T4 DNA  
polymerase (Biolabs) and the cDNA was digested with EcoRI. Since methylated dCTP was used during cDNA synthesis, the EcoRI site present in the tag was the only hemi-methylated site, hence the only site susceptible to EcoRI digestion. The cDNA was then size fractionated using exclusion chromatography  
20 (AcA, Biosepra) and fractions corresponding to cDNAs of more than 150 bp were pooled and ethanol precipitated. The cDNA was directionally cloned into the SmaI and EcoRI ends of the phagemid pBlueScript vector (Stratagene). The ligation mixture was electroporated into bacteria and propagated under appropriate antibiotic selection.

25

**EXAMPLE 5****Selection of Clones Having the Oligonucleotide Tag Attached Thereto**

Clones containing the oligonucleotide tag attached were then selected as follows. The plasmid DNAs containing 5' EST libraries made as described above were purified (Qiagen). A positive selection of the tagged clones was performed as follows. Briefly, in this selection procedure, the plasmid DNA  
30 was converted to single stranded DNA using gene II endonuclease of the phage F1 in combination with an exonuclease (Chang *et al.*, *Gene* 127:95-8, 1993) such as exonuclease III or T7 gene 6 exonuclease. The resulting single stranded DNA was then purified using paramagnetic beads as described by Fry *et al.*, *Biotechniques*, 13: 124-131, 1992. In this procedure, the single stranded DNA was hybridized with a biotinylated oligonucleotide having a sequence corresponding to the 3' end of the oligonucleotide tag.  
35 Clones including a sequence complementary to the biotinylated oligonucleotide were captured by incubation with streptavidin coated magnetic beads followed by magnetic selection. After capture of the positive clones, the plasmid DNA was released from the magnetic beads and converted into double

stranded DNA using a DNA polymerase such as the Thermosequenase obtained from Amersham Pharmacia Biotech. The double stranded DNA was then electroporated into bacteria. The percentage of positive clones having the 5' tag oligonucleotide was estimated using dot blot analysis to typically be between 90 and 98%.

- 5        Following electroporation, the libraries were ordered in 384-microtiter plates (MTP). A copy of the MTP was stored for future needs. Then the libraries were transferred into 96 MTP and sequenced as described below.

### EXAMPLE 6

#### 10        Sequencing of Inserts in Selected Clones

Plasmid inserts were first amplified by PCR on PE-9600 thermocyclers (Perkin-Elmer, Applied Biosystems Division, Foster City, CA), using standard SETA-A and SETA-B primers (Genset SA), AmpliTaqGold (Perkin-Elmer), dNTPs (Boehringer), buffer and cycling conditions as recommended by the Perkin-Elmer Corporation.

- 15        PCR products were then sequenced using automatic ABI Prism 377 sequencers (Perkin Elmer). Sequencing reactions were performed using PE 9600 thermocyclers with standard dye-primer chemistry and ThermoSequenase (Amersham Pharmacia Biotech). The primers used were either T7 or 21M13 (available from Genset SA) as appropriate. The primers were labeled with the JOE, FAM, ROX and TAMRA dyes. The dNTPs and ddNTPs used in the sequencing reactions were purchased from  
20   Boehringer. Sequencing buffer, reagent concentrations and cycling conditions were as recommended by Amersham.

- Following the sequencing reaction, the samples were precipitated with ethanol, resuspended in formamide loading buffer, and loaded on a standard 4% acrylamide gel. Electrophoresis was performed for 2.5 hours at 3000V on an ABI 377 sequencer, and the sequence data were collected and analyzed  
25   using the ABI Prism DNA Sequencing Analysis Software, version 2.1.2.

### EXAMPLE 7

#### Obtaining 5' ESTs from Extended cDNA libraries

##### Obtained from mRNA with Intact 5' Ends

- 30        Alternatively, 5'ESTs may be isolated from other cDNA or genomic DNA libraries. Such cDNA or genomic DNA libraries may be obtained from a commercial source or made using other techniques familiar to those skilled in the art. One example of such cDNA library construction, a full-length cDNA library, is as follows.

- PolyA+ RNAs are prepared and their quality checked as described in Example 1. Then, the  
35   caps at the 5' ends of the polyA+ RNAs are specifically joined to an oligonucleotide tag as described in Example 2. The oligonucleotide tag may contain a restriction site such as Eco RI to facilitate further

subcloning procedures. Northern blotting is then performed to check the size of mRNAs having the oligonucleotide tag attached thereto and to ensure that the mRNAs are actually tagged.

First strand synthesis is subsequently carried out for mRNAs joined to the oligonucleotide tag as described in Example 3 above except that the random nonamers are replaced by an oligo-dT primer. For instance, this oligo-dT primer may contain an internal tag of 4 nucleotides which is different from one tissue to the other. Following second strand synthesis using a primer contained in the oligonucleotide tag attached to the 5' end of mRNA, the blunt ends of the obtained double stranded full-length DNAs are modified into cohesive ends to facilitate subcloning. For example, the extremities of full-length cDNAs may be modified to allow subcloning into the Eco RI and Hind III sites of a Bluescript vector using the Eco RI site of the oligonucleotide tag and the addition of a Hind III adaptor to the 3' end of full-length cDNAs.

The full-length cDNAs are then separated into several fractions according to their sizes using techniques familiar to those skilled in the art. For example, electrophoretic separation may be applied in order to yield 3 or 6 different fractions. Following gel extraction and purification, the cDNA fractions are subcloned into appropriate vectors, such as Bluescript vectors, transformed into competent bacteria and propagated under appropriate antibiotic conditions. Subsequently, plasmids containing tagged full-length cDNAs are positively selected as described in Example 5.

The 5' end of full-length cDNAs isolated from such cDNA libraries may then be sequenced as described in Example 6 to yield 5'ESTs.

## II. Computer Analysis of the Isolated 5' ESTs: Construction of the SignalTag™ Database

The sequence data from the cDNA libraries made as described above were transferred to a database, where quality control and validation steps were performed. A base-caller, working using a Unix system, automatically flagged suspect peaks, taking into account the shape of the peaks, the inter-peak resolution, and the noise level. The base-caller also performed an automatic trimming. Any stretch of 25 or fewer bases having more than 4 suspect peaks was considered unreliable and was discarded. Sequences corresponding to cloning vector or ligation oligonucleotides were automatically removed from the EST sequences. However, the resulting EST sequences may contain 1 to 5 bases belonging to the above mentioned sequences at their 5' end. If needed, these can easily be removed on a case to case basis.

Following sequencing as described above, the sequences of the 5' ESTs were entered in a database for storage and manipulation as described below. Before searching the ESTs in the database for sequences of interest, ESTs derived from mRNAs which were not of interest were identified. Briefly, such undesired sequences may be of three types. First, contaminants of either endogenous (ribosomal RNAs, transfer RNAs, mitochondrial RNAs) or exogenous (prokaryotic RNAs and fungal RNAs) origins were identified. Second, uninformative sequences, namely redundant sequences, small sequences and highly degenerate sequences were identified. Third, repeated sequences (Alu, L1, THE

and MER repeats, SSTR sequences or satellite, micro-satellite, or telomeric repeats) were identified and masked in further processing.

In order to determine the accuracy of the sequencing procedure as well as the efficiency of the 5' selection described above, the analyses described in Examples 8 and 9 respectively were performed on 5' ESTs obtained from the database following the elimination of endogenous and exogenous contaminants and following the masking of repeats.

### EXAMPLE 8

#### Measurement of Sequencing Accuracy by Comparison to Known Sequences

To further determine the accuracy of the sequencing procedure described in Example 6, the sequences of 5' ESTs derived from known sequences were identified and compared to the original known sequences. First, a FASTA analysis with overhangs shorter than 5 bp on both ends was conducted on the 5' ESTs to identify those matching an entry in the public human mRNA database available at the time of filing the priority applications. The 5' ESTs which matched a known human mRNA were then realigned with their cognate mRNA and dynamic programming was used to include substitutions, insertions, and deletions in the list of "errors" which would be recognized. Errors occurring in the last 10 bases of the 5' EST sequences were ignored to avoid the inclusion of spurious cloning sites in the analysis of sequencing accuracy. This analysis revealed that the sequences incorporated in the database had an accuracy of more than 99.5%.

20

### EXAMPLE 9

#### Determination of Efficiency of 5' EST Selection

To determine the efficiency at which the above selection procedures isolated 5' ESTs which included sequences close to the 5' end of the mRNAs from which they derived, the sequences of the ends of the 5' ESTs derived from the elongation factor 1 subunit  $\alpha$  and ferritin heavy chain genes were compared to the known cDNA sequences of these genes. Since the transcription start sites of both genes are well characterized, they may be used to determine the percentage of derived 5' ESTs which included the authentic transcription start sites. For both genes, more than 95% of the obtained 5' ESTs actually included sequences close to or upstream of the 5' end of the corresponding mRNAs.

To extend the analysis of the reliability of the procedures for isolating 5' ESTs from ESTs in the database, a similar analysis was conducted using a database composed of human mRNA sequences extracted from GenBank database release 97 for comparison. The 5' ends of more than 85% of 5' ESTs derived from mRNAs included in the GenBank database were located close to the 5' ends of the known sequence. As some of the mRNA sequences available in the GenBank database are deduced from genomic sequences, a 5' end matching with these sequences will be counted as an internal match. Thus, the method used here underestimates the yield of ESTs including the authentic 5' ends of their corresponding mRNAs.

**EXAMPLE 10**Calculation of Novelty Indices for 5'EST Libraries

In order to evaluate the novelty of 5'EST libraries, the following analysis was performed. For  
5 each sequenced 5'EST library, the sequences were clustered by the 5' end. Each sequence in the library  
was compared to the others and the longest sequence found in the cluster was used as representative of  
the group. A novelty rate (NR) was then defined as:  $NR = 100 \times (\text{Number of new unique sequences found in the library} / \text{Total number of sequences from the library})$ . Typically, novelty rating ranged  
between 10% and 41% depending on the tissue from which the 5'EST library was obtained. For most of  
10 the libraries, the random sequencing of 5' EST libraries was pursued until the novelty rate reached 20%.

**EXAMPLE 11**Generation of Consensus Contigated 5' ESTs

Since the cDNA libraries made above include multiple 5' ESTs derived from the same mRNA,  
15 overlapping 5'ESTs may be assembled into continuous sequences. The following method describes how  
to efficiently align multiple 5'ESTs in order to yield not only consensus contigated 5'EST sequences for  
mRNAs derived from different genes but also consensus contigated 5'EST sequences for different  
mRNAs, so called variants, transcribed from the same gene such as alternatively spliced mRNAs.

The whole set of sequences was first partitioned into small clusters containing sequences  
20 which exhibited perfect matches with each other on a given length and which derived from a small  
number of different genes. Some 5'EST sequences, so called singletons, were not aligned using this  
approach because they were not homologous to any other sequence.

Thereafter, all variants of a given gene were identified in each cluster using a proprietary  
software. 5'EST sequences belonging to the same variant were then contigated and consensus  
25 contigated 5'EST sequences generated for each variant. All consensus contigated 5' EST sequences  
were subsequently compared to the whole set of individual 5'EST sequences used to obtained them.

If desired, the consensus contigated 5'EST sequences may be verified by identifying clones  
in nucleic acid samples derived from biological tissues, such as cDNA libraries, which hybridize to  
the probes based on the sequences of the consensus contigated 5'ESTs using any methods described  
30 herein and sequencing those clones.

Application of this alignment method to a selected set of 5'ESTs free from endogenous  
contaminants and uninformative sequences, and following the masking of repeats, yielded consensus  
contigated 5'EST sequences or variants of clustered genes encompassing many individual 5'ESTs.  
Both non aligned 5'ESTs, *i.e.* singletons, and consensus contigated 5'ESTs were then compared to  
35 already known sequences and those sequences matching human mRNA sequences were eliminated  
from further analysis.

## EXAMPLE 12

Identification of Open Reading Frames in 5' ESTs

Subsequently, consensus contigated 5'ESTs and 5'ESTs were screened to identify those having an open reading frame (ORF).

5 Such open reading frames were simply defined as uninterrupted nucleic acid sequences longer than 45 nucleotides and beginning with an ATG codon.

Alternatively, the nucleic acid sequence was first divided into several subsequences which coding propensity was evaluated separately using one or several different methods known to those skilled in the art such as the evaluation of N-mer frequency and its variants (Fickett and Tung, 10 *Nucleic Acids Res*;20:6441-50 (1992)) or the Average Mutual Information method (Grosse *et al*, International Conference on Intelligent Systems for Molecular Biology, Montreal, Canada. June 28-July 1, 1998). Each of the scores obtained by the techniques described above were then normalized by their distribution extremities and then fused using a neural network into a unique score that represents the coding probability of a given subsequence. The coding probability scores obtained for 15 each subsequence, thus the probability score profiles obtained for each reading frame, was then linked to the initiation codons present on the sequence. For each open reading frame, defined as a nucleic acid sequence beginning with an ATG codon, an ORF score was determined. Preferably, this score is the sum of the probability scores computed for each subsequence corresponding to the considered ORF in the correct reading frame corrected by a function that negatively accounts for 20 locally high score values and positively accounts for sustained high score values. The most probable ORF with the highest score was selected.

In some embodiments, nucleic acid sequences encoding an "incomplete ORF", as referred therein, namely an open reading frame in which a start codon has been identified but no stop codon has been identified, were obtained.

25 In other embodiments, nucleic acid sequences encoding a "complete ORF", as used therein, namely an open reading frame in which a start codon and a stop codon have been identified, are obtained.

In a preferred embodiment, open reading frames encoding polypeptides of at least 50 amino acids were obtained.

30 To confirm that the chosen ORF actually encodes a polypeptide, the consensus contigated 5'EST or 5'EST may be used to obtain an extended cDNA using any of the techniques described therein, and especially those described in Examples 19 and 20. Then, such obtained extended cDNAs may be screened for the most probable open reading frame using any of the techniques described therein. The amino acid sequence of the ORF encoded by the consensus contigated 5'EST or 5'EST may then be 35 compared to the amino acid sequence of the ORF encoded by the extended cDNA using any of the algorithms and parameters described therein in order to determine whether the ORF encoded by the extended cDNA is basically the same as the one encoded by the consensus contigated 5'EST or 5'EST.

Alternatively, to confirm that the chosen ORF actually encodes a polypeptide, the consensus contiguated 5'EST or 5'EST may be used to obtain an extended cDNA using any of the techniques described therein, and especially those described in Examples 19 and 20. Such an extended cDNA may then be inserted into an appropriate expression vector and used to express the polypeptide encoded by the extended cDNA as described therein. The expressed polypeptide may be isolated, purified, or enriched as described therein. Several methods known to those skilled in the art may then be used to determine whether the expressed polypeptide is the one actually encoded by the chosen ORF, therein referred to as the expected polypeptide. Such methods are based on the determination of predictable features of the expressed polypeptide, including but not limited to its amino acid sequence, its size or its charge, and the comparison of these features to those predicted for the expected polypeptide. The following paragraphs present examples of such methods.

One of these methods consists in the determination of at least a portion of the amino acid sequence of the expressed polypeptide using any technique known to those skilled in the art. For example, the amino-terminal residues may be determined using techniques either based on Sanger's technique of acid hydrolysis of a polypeptide which N-terminal residue has been covalently labeled or using techniques based on Edman degradation of polypeptides which N-terminal residues are sequentially labeled and cleaved from the polypeptide of interest. The amino acid sequence of the expressed polypeptide may then be compared to the one predicted for the expected polypeptide using any algorithm and parameters described therein.

Alternatively, the size of the expressed polypeptides may be determined using techniques familiar to those skilled in the art such as Coomassie blue or silver staining and subsequently compared to the size predicted for the expected polypeptide. Generally, the band corresponding to the expressed polypeptide will have a mobility near that expected based on the number of amino acids in the open reading frame of the extended cDNA. However, the band may have a mobility different than that expected as a result of modifications such as glycosylation, ubiquitination, or enzymatic cleavage.

Alternatively, specific antibodies or antipeptides may be generated against the expected polypeptide as described in Example 34 and used to perform immunoblotting or immunoprecipitation studies against the expressed polypeptide. The presence of a band in samples from cells containing the expression vector with the extended cDNA which is absent in samples from cells containing the expression vector encoding an irrelevant polypeptide indicates that the expected polypeptide or portion thereof is being expressed. Generally, the band corresponding to the expressed polypeptide will have a mobility near that expected based on the number of amino acids in the open reading frame of the extended cDNA. However, the band may have a mobility different than that expected as a result of modifications such as glycosylation, ubiquitination, or enzymatic cleavage

35

### EXAMPLE 13

#### Identification of Potential Signal Sequences in 5' ESTs

The 5'ESTs or consensus contigated 5'ESTs found to encode an ORF were then searched to identify potential signal motifs using slight modifications of the procedures disclosed in Von Heijne, *Nucleic Acids Res.* 14:4683-4690, 1986. Those sequences encoding a 15 amino acid long stretch with a score of at least 3.5 in the Von Heijne signal peptide identification matrix were considered to possess a signal sequence. Those nucleic acid sequences which match a known human mRNA or EST sequence and have a 5' end located downstream of the known 5' end, preferably by more than 20 nucleotides, were excluded from further analysis. The remaining nucleic acids having signal sequences therein were included in a database called SignalTag™.

10

#### EXAMPLE 14

##### Confirmation of Accuracy of Identification of Potential Signal Sequences in 5' ESTs

The accuracy of the above procedure for identifying signal sequences encoding signal peptides was evaluated by applying the method to the 43 amino acids located at the N terminus of all human SwissProt proteins. The computed Von Heijne score for each protein was compared with the known characterization of the protein as being a secreted protein or a non-secreted protein. In this manner, the number of non-secreted proteins having a score higher than 3.5 (false positives) and the number of secreted proteins having a score lower than 3.5 (false negatives) could be calculated.

Using the results of the above analysis, the probability that a peptide encoded by the 5' region of the mRNA is in fact a genuine signal peptide based on its Von Heijne's score was calculated based on either the assumption that 10% of human proteins are secreted or the assumption that 20% of human proteins are secreted. The results of this analysis are shown in Figure 2.

Using the above method of identification of secretory proteins, 5' ESTs of the following polypeptides known to be secreted were obtained: human glucagon, gamma interferon induced monokine precursor, secreted cyclophilin-like protein, human pleiotropin, and human biotinidase precursor. Thus, the above method successfully identified those 5' ESTs which encode a signal peptide.

To confirm that the signal peptide encoded by the 5' ESTs or consensus contigated 5' ESTs actually functions as a signal peptide, the signal sequences from the 5' ESTs or consensus contigated 5' ESTs may be cloned into a vector designed for the identification of signal peptides. Such vectors are designed to confer the ability to grow in selective medium only to host cells containing a vector with an operably linked signal sequence. For example, to confirm that a 5' EST or consensus contigated 5' EST encodes a genuine signal peptide, the signal sequence of the 5' EST or consensus contigated 5' EST may be inserted upstream and in frame with a non-secreted form of the yeast invertase gene in signal peptide selection vectors such as those described in U.S. Patent No. 5,536,637. Growth of host cells containing signal sequence selection vectors with the correctly inserted 5' EST or consensus contigated 5' EST signal sequence confirms that the 5' EST or consensus contigated 5' ESTs encodes a genuine signal peptide.



Alternatively, the presence of a signal peptide may be confirmed by cloning the extended cDNAs obtained using the ESTs or consensus contigated 5' ESTs into expression vectors such as pXT1 as described below, or by constructing promoter-signal sequence-reporter gene vectors which encode fusion proteins between the signal peptide and an assayable reporter protein. After introduction of these vectors into a suitable host cell, such as COS cells or NIH 3T3 cells, the growth medium may be harvested and analyzed for the presence of the secreted protein. The medium from these cells is compared to the medium from control cells containing vectors lacking the signal sequence or extended cDNA insert to identify vectors which encode a functional signal peptide or an authentic secreted protein.

## EXAMPLE 15

### Analysis of the Sequences of the Invention

The set of the nucleic acid sequences of the invention (SEQ ID NOs. 24-811 and 1600-1622) was obtained as described in Example 11. Subsequently, the most probable open reading frame was determined and signal sequences were searched, as described in Examples 12 and 13, for all sequences of the invention.

The nucleotide sequences of the SEQ ID NOs. 24-811 and 1600-1622 and the polypeptides sequences encoded by SEQ ID NOs. 24-811 (*i.e.* polypeptide sequences of SEQ ID NOs. 812-1599) are provided in the appended sequence listing which structure is as follows.

SEQ ID NOs. 24-728 are nucleic acids having an incomplete ORF which encodes a signal peptide. The locations of the incomplete ORFs and sequences encoding signal peptides are listed in the accompanying Sequence Listing. In addition, the von Heijne score of the signal peptide computed as described in Example 13 is listed as the "score" in the accompanying Sequence Listing. The sequence of the signal-peptide is listed as "seq" in the accompanying Sequence Listing. The "/" in the signal peptide sequence indicates the location where proteolytic cleavage of the signal peptide occurs to generate a mature protein.

SEQ ID NOs. 729-765 are nucleic acids having an incomplete ORF in which no sequence encoding a signal peptide has been identified to date. However, it remains possible that subsequent analysis will identify a sequence encoding a signal peptide in these nucleic acids. The locations of the incomplete ORFs are listed in the accompanying Sequence Listing.

SEQ ID NOs. 766-792 are nucleic acids having a complete ORF which encodes a signal peptide. The locations of the complete ORFs and of the signal peptides, the von Heijne score of the signal peptide, the sequence of the signal-peptide and the proteolytic cleavage site are indicated as described above.

SEQ ID NOs. 793-811 are nucleic acids having a complete ORF in which no sequence encoding a signal peptide has been identified to date. However, it remains possible that subsequent analysis will

identify a sequence encoding a signal peptide in these nucleic acids. The locations of the complete ORFs are listed in the accompanying Sequence Listing.

SEQ ID NOs. 812-1516 are "incomplete polypeptide sequences" which include a signal peptide. "Incomplete polypeptide sequences" are polypeptide sequences encoded by nucleic acids in which a start codon has been identified but no stop codon has been identified. These polypeptides are encoded by the nucleic acids of SEQ ID NOs. 24-728. The location of the signal peptide, the von Heijne score of the signal peptide, the sequence of the signal-peptide and the proteolytic cleavage site are indicated as described above.

SEQ ID NOs. 1517-1553 are incomplete polypeptide sequences in which no signal peptide has been identified to date. However, it remains possible that subsequent analysis will identify a signal peptide in these polypeptides. These polypeptides are encoded by the nucleic acids of SEQ ID NOs. 729-765.

SEQ ID NOs. 1554-1580 are "complete polypeptide sequences" which include a signal peptide. "Complete polypeptide sequences" are polypeptide sequences encoded by nucleic acids in which a start codon and a stop codon have been identified. These polypeptides are encoded by the nucleic acids of SEQ ID NOs. 766-792. The location of the signal peptide, the von Heijne score of the signal peptide, the sequence of the signal-peptide and the proteolytic cleavage site are indicated as described above..

SEQ ID NOs. 1581-1599 are complete polypeptide sequences in which no signal peptide has been identified to date. However, it remains possible that subsequent analysis will identify a signal peptide in these polypeptides. These polypeptides are encoded by the nucleic acids of SEQ ID NOs. 793-811.

SEQ ID NOs. 1600-1622 are nucleic acid sequences in which no open reading frame has been conclusively identified to date. However, it remains possible subsequent analysis will identify an open reading frame in these nucleic acids.

In the accompanying Sequence Listing, all instances of the symbol "n" in the nucleic acid sequences mean that the nucleotide can be adenine, guanine, cytosine or thymine. In some instances the polypeptide sequences in the Sequence Listing contain the symbol "Xaa." These "Xaa" symbols indicate either (1) a residue which cannot be identified because of nucleotide sequence ambiguity or (2) a stop codon in the determined sequence where applicants believe one should not exist (if the sequence were determined more accurately). In some instances, several possible identities of the unknown amino acids may be suggested by the genetic code.

In the case of secreted proteins, it should be noted that, in accordance with the regulations governing Sequence Listings, in the appended Sequence Listing, the full protein (*i.e.* the protein containing the signal peptide and the mature protein) extends from an amino acid residue having a negative number through a positively numbered C-terminal amino acid residue. Thus, the first amino acid of the mature protein resulting from cleavage of the signal peptide is designated as amino acid

number 1, and the first amino acid of the signal peptide is designated with the appropriate negative number.

If one of the nucleic acid sequences of SEQ ID NOs. 24-811 and 1600-1622 are suspected of containing one or more incorrect or ambiguous nucleotides, the ambiguities can readily be resolved by resequencing a fragment containing the nucleotides to be evaluated. If one or more incorrect or ambiguous nucleotides are detected, the corrected sequences should be included in the clusters from which the sequences were isolated, and used to compute other consensus contigated sequences on which other ORFs would be identified. Nucleic acid fragments for resolving sequencing errors or ambiguities may be obtained from deposited clones or can be isolated using the techniques described herein.

Resolution of any such ambiguities or errors may be facilitated by using primers which hybridize to sequences located close to the ambiguous or erroneous sequences. For example, the primers may hybridize to sequences within 50-75 bases of the ambiguity or error. Upon resolution of an error or ambiguity, the corresponding corrections can be made in the protein sequences encoded by the DNA containing the error or ambiguity. The amino acid sequence of the protein encoded by a particular clone can also be determined by expression of the clone in a suitable host cell, collecting the protein, and determining its sequence.

In addition, if one of the sequences of SEQ ID NOs. 812-1599 is suspected of containing a truncated ORF as the result of a frameshift in the sequence, such frameshifting errors may be corrected by combining the following two approaches. The first one involves thorough examination of all double predictions, *i.e.* all cases where the probability scores for two ORFs located on different reading frames are high and close, preferably different by less than 0.4. The fine examination of the region where the two possible ORFs overlap may help to detect the frameshift. In the second approach, homologies with known proteins are used to correct suspected frameshifts.

Of the identified clusters, some were shown to be multivariant, *i.e.* to contain several variants of the same gene. Table I gives for each of the multivariant clusters named by its internal reference (first column), the list of all variant consensus contigated 5'ESTs (second column), each being represented by a different sequence identification number.

TABLE I

Cluster Internal Reference	SEQ ID NOs of Variants
C1	687, 791
C2	744, 798
C3	640, 811
C4	59, 66
C5	84, 97

C6	287, 289
C7	286, 775, 777
C8	762, 768
C9	783, 784
C10	80, 1603
C11	655, 736
C12	805, 806

Table II provides a list preferred polynucleotide fragments which are derivatives of the consensus contigated 5'ESTs. As used herein the term "polynucleotide described in Table II" refers to the all of the preferred polynucleotide fragments defined in Table II in the following manner. The fragments are referred to by their SEQ ID numbers in the first column. The preferred polynucleotide fragments are then defined by a range of nucleotide positions from the SEQ IDs of the consensus contigated 5'ESTs as indicated in the second column entitled "positions of preferred fragments." The preferred polynucleotide fragments correspond to the individual 5'ESTs aligned to obtain the consensus contigated 5'EST and to those filed in the priority documents. The third column entitled "variant nucleotides" describes the nucleotide sequence variations observed between the consensus contigated 5'EST and preferred nucleic acid fragments as follows:

A) Substitutions in the sequence of a consensus contigated 5'EST to derive a preferred polynucleotide fragment are denoted by an "S", followed by a number indicating the first nucleotide position in a specific SEQ ID to be substituted in a string of substituted nucleotides or the position of the substituted nucleotide in the case of a single substituted nucleotide. Then there is a coma followed by one or more lower case letters indicating the identity of the nucleotide(s) occurring in the substituted position(s). For example, SEQ ID NO: 3401; Position of preferred fragments: 1-250; Variant nucleotides S45,atc would indicate that a preferred polynucleotide fragment had the sequence of positions 1 to 250 of SEQ ID NO. 3401, except that the nucleotides at positions 45, 46, and 47 were substituted with A, T, and C, respectively, in the preferred polynucleotide as compared with the sequence of SEQ ID No. 3401.

B) Insertions in the sequence of a consensus contigated 5'EST to derive a preferred polynucleotide fragment are denoted by an "I", followed by a number indicating the nucleotide position in a specific SEQ ID after which a string of nucleotides is inserted or the position after which the nucleotide is inserted in the case of a single inserted nucleotide. Then there is a coma followed by one or more lower case letters indicating the identity of the nucleotide(s) occurring in the inserted position(s). For example, SEQ ID NO: 7934; Position of preferred fragments: 1-500; Variant nucleotides I36,gataca would indicate that a preferred polynucleotide fragment had the sequence of positions 1 to 500 of SEQ ID NO. 7934, except that after the nucleotides at position 36 a GATACA string of nucleotides is inserted in the preferred polynucleotide as compared with the sequence of SEQ ID No. 7934.

C) Deletions in the sequence of a consensus contigated 5'EST to derive a preferred nucleic acid fragment are denoted by an "D", followed by a number indicating the first nucleotide position in a specific SEQ ID to be deleted in a string of deleted nucleotides or the position of the deleted nucleotide in the case of a single deleted nucleotide. Then there is a comma followed by number indicating the number of nucleotide(s) deleted from the sequence provided in the sequence ID. For example, SEQ ID NO: 5398; Position of preferred fragments: 56-780; Variant nucleotides D114,5 would indicate that a preferred polynucleotide fragment had the sequence of positions 56 to 780 of SEQ ID NO. 5398, except that the nucleotides in positions 114 to 118 had been deleted in the preferred polynucleotide as compared with the sequence of SEQ ID No. 5398.

The present invention encompasses isolated, purified, or recombinant nucleic acids which consist of, consist essentially of, or comprise a contiguous span of at least 8, 10, 12, 15, 18, 20, 25, 35, 40, 50, 70, 80, 100, 250, or 500 nucleotides in length, to the extent that a contiguous span of these lengths is consistent with the lengths of the particular polynucleotide, of a polynucleotide described in Table II, or a sequence complementary thereto, wherein said polynucleotide described in Table II is selected individually or in any combination from the polynucleotides described in Table II. The present invention also encompasses isolated, purified, or recombinant nucleic acids which consist of or consist essentially of a polynucleotide described in Table II, or a sequence complementary thereto, wherein said polynucleotide is selected individually or in any combination from the polynucleotides described in Table II. The present invention further encompasses isolated or purified polypeptides which consist of, consist essentially of, or comprise a contiguous span of at least 8, 10, 12, 15, 18, 20, 25, 35, 40, 50, 70, 80, or 100 amino acids encoded by a polynucleotide described in Table II.

Table II

SEQ ID NO.	Positions of Preferred Fragments	Variant nucleotides
35	1-423	S124, s; I135, a; S293, w; I363, a; S377, r; D424, 15
41	1-427	I117, m; S120, r; S124, g; D373, l; S376, b; S378, b; I427, gggg; D428, 109
43	1-276	S114, m; S118, rg; S123, r; S139, nr; I142, t; D148, l; D152, l; I228, t; I276, gg; D277, 136
45	126-420	D1, 125; I420, ggg; D421, 100
46	1-255	S139, r; I145, r; S146, mm; S150, ar; S254, g; D256, 128
48	4-437	D1, 3; S49, a; S55, g; S79, a; S90, a; I437, tctctg
59	1-471	S26, a; S44, t; S48, t; S109, a; S191, t; S200, gc; S203, a; S210, g; S237, a; S240, g; S255, a; S272, a; S277, a; S279, a; S284, t; S297, g; S305, g; S316, a; I471, ggtca
66	1-428	I428, tactgggg

82	1-399	S251, t; S277, d; I399, aagccggg
84	5-488	D1, 4; S210, g; S293, a; S325, g; S339, a; S348, g; S353, g; S395, g; I488, cacca
93	1-508	I508, gattt
96	26-315	D1, 25; S28, a; S62, c; I315, cagatgg
97	4-460	D1, 3; S19, g; S31, g; S114, gt; S118, a; S123, tc; S127, c; S132, a; S186, g; S190, c; S203, t; S210, g; S232, c; I460, acgtt
105	1-281	S273, a; I281, g; D282, 211
114	10-315	I0, t; D1, 9; S91, m; S267, n; S276, w; S292, h; S295, m; I315, tggg; D316, 19
118	1-145	S57, d; S126, d; I145, ccctc
120	2-348	D1, 1; S104, t; I348, g; D349, 38
121	1-190	I121, c; I190, ccctt
123	1-353	I117, m; I186, w; S187, y; I353, caccgggg
124	1-249	I249, ggrvgggg
125	114-375	D1, 113; S206, wn; I231, a; I375, ccctagg
126	1-437	S297, cc; S307, tg; S312, a; S318, g; S341, a; S351, t; S353, g; S383, c; S387, a; D404, 1
136	82-428	D1, 81; I428, aaagtg
139	1-268	I268, gggaaggg
148	6-405	D1, 5; I405, ggtgt
159	1-230	S227, ta; I230, ccctggg
165	3-256	I0, tat; D1, 2; I17, c; S18, t; S111, d; I115, t; S123, r; I256, aagccggg
170	1-280	I103, t; S104, c; I111, t; I280, cgttcggg
194	1-215	S50, s; S186, sn; S199, k; I215, gcagcggg
213	1-158	S128, m; I132, w; S143, d; I158, tgcccggg
223	3-431	D1, 2; S28, s; S79, c; S82, s; S308, nr; S328, nb; I431, ccggc
247	1-359	I76, gttt; I359, tccttg
258	1-236	S72, r; S81, g; S197, s; I205, ss; S232, k; I236, acttcggg
264	5-283	D1, 4; S64, g; S122, m; S134, yy; I137, c; I151, t; I283, gttgc
269	1-143	S111, s; I143, ggggcggg
286	5-207	D1, 4; S204, a; S206, c; I207, gg; D208, 567
287	1-277	S114, r; I125, t; S131, ag; S256, tg; S259, tt; S262, at; S267, t; S269, c; S273, c; I277, ccggg; D278, 337
289	69-416	D1, 68; I416, agccaggg
289	1-278	S114, r; I125, t; S131, ag; S277, c; I278, cggg; D279, 138
292	20-254	D1, 19; I254, aaagagg
293	1-414	I414, tagcag
300	1-285	S16, m; S67, y; I285, baccacggg; D286, 1
349	23-431	D1, 22; I118, a; S214, y; I431, caactgg
350	3-386	D1, 2; S42, w; I263, c; I386, gggat
368	3-446	D1, 2; I446, tctct
385	1-193	I35, t; I108, t; I134, r; S135, a; S137, r; S143, w; I178, c; I193, gagcgggg
411	6-391	D1, 5; S17, r; S27, t; S334, y; D392, 244
412	1-185	S49, s; S127, s; I185, gctggg; D186, 150

415	2-229	D1, l; S3, a; I229, caaatggg
435	1-386	S4, s; I386, ccggg
436	4-472	D1, 3; S61, sa; D238, l; S239, s; I472, agtgtgg
437	1-340	I340, ggg; D341, 129
441	1-409	S109, smag; I409, cgcacggg
454	1-492	S72, nn; S115, t; S121, bwy; S181, yn; I492, gagtc
455	1-177	I14, w; I16, a; I177, gagctggg
459	1-311	S39, n; S74, rg; I311, accatggg
460	1-425	I425, agtac
461	5-420	D1, 4; I420, tcgtc
481	1-429	I10, w; S262, d; S333, n; I429, ctccaggg
489	1-414	D72, l; S117, n; S396, d; I414, ggaca
496	1-215	I215, ttttcggg
501	1-430	S275, n; I430, aggat
502	91-413	D1, 90; I413, aaacgggg
504	21-420	D1, 20; S47, w; S83, n; I280, n; S281, na; S292, v; S314, sm; S368, ww; S373, w; I420, cccca
505	18-457	D1, 17; D36, l; S182, g; S273, n; S283, a; S416, bh; I457, ctcga
514	1-303	I303, accca
515	1-455	S11, t; I12, n; S30, r; S256, wr; I333, t; I455, cataa
517	24-453	D1, 23; I453, agagcggg
519	1-275	I119, gt; S125, w; I129, w; S133, k; S137, k; S167, k; I275, gcccc
522	1-313	I313, agcgtggg
526	4-366	I0, t; D1, 3; I366, ggcccggg
530	1-434	S328, g; I434, aagat
535	1-379	S128, g; S162, m; D380, 5
561	2-341	D1, l; I341, raagagg
568	1-246	I118, g; S137, g; I246, aaaccggg
570	1-207	I207, tttt
576	1-288	I34, c; I288, cccgtgg
588	1-390	S218, a; S224, k; S314, dh; S358, s; D376, l; I390, atg; D391, 23
597	31-274	D1, 30; S49, n; I274, tccatgg
606	1-354	I141, g; D174, l; S229, rr; D355, 72
627	1-415	S7, a; I415, cattt
634	1-178	D179, 212
640	6-428	D1, 5; D429, 79
641	64-483	D1, 63; I165, d; D183, l; S185, y; S253, t; D279, 2; S416, a; I483, atata
655	1-280	S58, c; I84, g; S88, k; S204, ac; S244, g; S247, g; I280, ggg; D281, 90
672	34-489	D1, 33; S316, k; S331, k; S333, w; S486, g; S488, c; D490, 4
687	116-473	D1, 115; S142, n; I473, cctcgggg
697	1-202	S142, s; S144, sr; S148, d; S152, d; I155, a; I164, a; S174, k; I202, gcc; D203, 291
708	8-384	D1, 7; S104, b; I384, gaaaa
710	1-167	S40, k; S49, db; I167, tatct

722	1-191	I125, c; I191, tttt
723	1-316	I316, aggg; D317, 157
729	15-373	D1, 14; S139, t; I373, cgcag; D374, 99
730	29-372	D1, 28; I155, g; S192, ka; S333, d; I372, m; D373, 93
731	1-290	S10, kk; S30, b; S32, t; S92, t; S197, dy; S278, g; I290, aggg; D291, 55
732	8-277	D1, 7; I113, a; S127, w; I131, s; S132, r; S156, w; S160, r; S211, n; S215, w; I247, a; D278, 121
733	20-375	D1, 19; S306, sbs; I325, h; S326, nr; S338, ywd; S344, v; I375, aggg; D376, 68
734	1-359	D66, 1; D360, 14
735	25-322	D1, 24; S30, r; I193, a; I322, ccaaggg
736	9-181	D1, 8; S58, g; I181, aactaggg
737	1-160	S97, ta; I160, aggtc
738	1-227	D228, 7
739	45-514	D1, 44; S178, s; I182, c; S436, dmn; S461, v; S476, c; S506, t; D515, 75
740	11-388	D1, 10; I388, cgacaggg
741	1-478	S118, s; S125, a; I126, s; S134, k; S421, vn; I478, aatsc
742	217-553	I0, tt; D1, 216; S286, r; S294, m; S311, r; S317, s; S338, r; S442, dm; S469, h; S476, r; S485, s; S491, w; I495, ht; S496, v; S513, r; D521, 1; S536, m; D554, 199
743	1-459	I11, s; S258, m; I270, m; I304, c; I308, amta; S313, c; S438, v; I459, agggag
744	25-316	D1, 24; S315, g; D317, 95
745	21-283	D1, 20; I40, g; S41, c; D123, 1; S181, sr; S227, r; I283, ccgcg; D284, 121
746	1-256	D257, 173
747	1-179	S134, w; S138, w; S140, kt; I179, cacca
748	1-235	S46, t; I72, t; S189, cc; S222, c; D236, 148
749	2-370	D1, 1; S32, cg; D144, 1; S341, g; D371, 76
750	18-410	I0, aag; D1, 17; I410, aatcc
751	22-355	D1, 21; D148, 1; S150, c; S152, a; S313, n; D356, 181
752	1-139	S50, t; I118, g; I139, ccct
753	1-189	S26, r; S115, s; I121, r; S122, r; S128, s; S143, r; I146, w; S156, r; D190, 4
754	1-395	S212, wd; I395, cggca
755	19-460	D1, 18; S26, c; S156, a; S253, n; I460, tagaagg
756	2-142	D1, 1; I106, gc; S107, t; S110, c; I142, ccaccggg
757	28-296	D1, 27; I119, s; I122, t; S128, s; S255, t; S267, m; D297, 66
758	11-368	D1, 10; I200, g; S201, c; S281, d; S317, c; I368, ccatcggg
759	19-452	D1, 18; S421, w; I452, a
760	25-175	D1, 24; S34, yk; I175, ccggg; D176, 120
761	1-212	I212, cactcggg
762	1-374	S320, s; S349, a; D375, 249
763	8-152	D1, 7; I152, acggg; D153, 109



764	1-160	I127, g; I145, g; I160, cgcccggg
765	137-313	D1, 136; S272, m; I279, s; S310, t; I313, ggg; D314, 203
766	1-320	S278, ag; S281, cagacc; S288, ta; S291, caag; S296, c; S317, m; I320, cggg; D321, 306
767	6-336	I0, aa; D1, 5; S149, w; S245, y; D337, 137
768	1-374	S320, s; D375, 299
769	53-435	D1, 52; S59, b; S344, nnkw; D436, 104
770	24-448	D1, 23; S25, g; S411, w; S416, m; D449, 31
771	1-370	S3, c; S180, m; S275, r; D371, 122
772	1-388	I299, c; S326, c; D389, 8
773	1-143	S18, c; S66, a; I143, ggg; D144, 274
774	1-347	S194, a; S205, c; I347, ggg; D348, 107
775	5-207	D1, 4; S111, tg; S158, g; S171, c; S191, a; S204, a; S206, c; I207, gg; D208, 324
776	1-368	I200, c; S201, a; S291, ta; I332, c
777	5-207	D1, 4; S204, a; S206, c; I207, gg; D208, 262
778	39-342	D1, 38; S184, r; D343, 126
779	4-360	D1, 3; S13, m; S15, c; S22, s; S24, m; S48, r; S56, s; S335, c; S345, rs; I360, ggg; D361, 119
780	1-472	I347, c; D473, 32
781	116-426	D1, 115; S219, m; S424, g; D427, 118
782	1-391	S386, k; D392, 64
783	1-453	D109, l; S110, y; S125, y; I128, g; S132, k; I453, ctctc
784	29-494	D1, 28; S72, r; D495, 93
785	99-461	D1, 98; S218, r; I461, gaccgggg
786	2-465	D1, 1; S8, y; S388, s; I398, g; S400, t; S403, at; S417, g; D466, 24
787	28-271	D1, 27; S99, t; S230, c; S266, ga; S269, c; I271, g; D272, 126
788	1-285	D280, 1; I285, g; D286, 310
789	1-209	S205, c; D210, 150
790	51-297	D1, 50; I297, ggggg; D298, 539
791	113-327	D1, 112; S218, g; I226, g; D280, 1; I327, cgcagg; D328, 224
792	17-218	D1, 16; S58, t; S217, t; I218, gggg; D219, 219
793	11-92	D1, 10; S91, c; I92, a; D93, 258
794	9-431	D1, 8; I431, taagt
795	30-341	D1, 29; I341, a; D342, 175
796	1-442	S17, w; S19, wr; D35, 1; S134, t; S264, n; S322, nr; S369, s; S420, s; S422, y; I442, tctcggg
797	1-420	S136, c; S150, c; I245, ccc; I420, ggagtg
798	25-316	D1, 24; S315, g; D317, 97
799	1-344	D345, 57
800	7-465	D1, 6; S59, k; S146, a; S186, km; I465, gttca
801	121-422	D1, 120; I269, c; S419, cc; I422, gg; D423, 207
802	46-477	D1, 45; S132, bn; I477, actac
803	15-467	D1, 14; S45, k; S65, t; S418, ys; D452, 1; D468, 119
804	1-341	S42, t; S97, d; S326, gtg; S331, tgt; S336, a;

		S338, c; I341, cccccggg; D342, 218
805	2-409	D1, 1; S334, d; I409, aggg; D410, 161
806	5-384	D1, 4; I384, actaa
807	1-301	S113, a; S117, c; S123, t; D128, 1; D134, 1; S282, g; S284, a; I301, gacggagggg; D302, 70
808	2-314	D1, 1; S306, g; I314, ggg; D315, 121
809	1-394	S53, g; S228, n; S272, vk; I301, g; I358, m; S368, nb; S375, w; I383, mm; I388, yt; I394, nhaccggg
810	6-205	I0, a; D1, 5; I141, t; I205, ggg; D206, 630
811	6-270	D1, 5; I270, gggg; D271, 115
1600	1-247	S45, m; S114, k; I122, m; S123, yc; S158, rr; S221, k; I247, ccccaggg
1601	1-225	S109, bm; S195, m; I225, tgcacggg
1602	23-245	D1, 22; D138, 1; S139, s; S242, t; S244, g; I245, g; D246, 13
1603	1-303	S71, c; D277, 1; I303, ggagggg; D304, 38
1604	1-242	S47, w; S50, c; S81, h; S85, d; S91, k; S106, r; I242, tgtggg; D243, 50
1605	2-225	D1, 1; S20, k; S91, c; I225, ggg; D226, 132
1606	15-293	D1, 14; S156, g; S193, g; I200, t; I293, acaaaggg
1607	1-361	S323, c; I361, cccca
1608	1-151	I151, taagggg; D152, 154
1609	1-242	S55, s; I135, a; S152, h; I242, cagtaggg
1610	1-196	I151, w; S190, k; I196, cctgtgg
1611	1-228	S115, k; S174, rk; I228, cgtttggg
1612	1-221	S108, v; I221, tgatcggg
1613	1-281	I66, w; I137, a; D282, 79
1614	1-171	S53, k; S76, k; I80, k; S81, kw; S86, r; S92, k; S126, k; I171, gccgagg
1615	2-193	D1, 1; S67, c; I121, s; S122, mm; S126, g; S130, r; S146, r; S156, gm; I193, cctca
1616	1-349	S251, ww; S259, rs; S275, k; I279, w; S285, y; S292, y; I320, m; I331, m; I338, w; I341, s; I349, accccggg
1617	1-129	I118, t; D130, 26
1618	1-184	D9, 1; D185, 1
1619	1-169	I122, t; I169, gcccgagg
1620	1-187	S106, k; S118, m; S122, cg; S132, k; D188, 59
1621	1-153	D125, 1; I131, ttt; S152, t; I153, gg; D154, 127
1622	1-400	S43, s; I126, g; I129, y; S353, d; I400, tatat

### EXAMPLE 16

#### Categorization of 5' ESTs and Consensus Contigated 5'ESTs

The nucleic acid sequences of the present invention (SEQ ID NOs. 24-811 and 1600-1622) were  
5 grouped based on their homology to known sequences as follows. All sequences were compared to  
EMBL release 57 and daily releases available at the time of filing using BLASTN. All matches with a  
minimum of 25 nucleotides with 90% homology were retrieved and used to compute Tables IV and V.

In some embodiments, 5'ESTs or consensus contigated 5'ESTs nucleic acid sequence do not match any known vertebrate sequence nor any publicly available EST sequence, thus being completely new.

In other embodiments, 5'ESTs or consensus contigated 5'ESTs match a known sequence.

- 5 Tables III and IV gives for each sequence of the invention in this category referred to by its sequence identification number in the first column, the positions of their preferred fragments in the second column entitled "Positions of preferred fragments." As used herein the term "polynucleotide described in Table III" refers to the all of the preferred polynucleotide fragments defined in Table III in this manner, and the term "polynucleotide described in Table IV" refers to the all of the preferred polynucleotides fragments
- 10 defined in Table IV in this manner. The present invention encompasses isolated, purified, or recombinant nucleic acids which consist of, consist essentially of, or comprise a contiguous span of at least 8, 10, 12, 15, 18, 20, 25, 35, 40, 50, 70, 80, 100, 250, or 500 nucleotides in length, to the extent that a contiguous span of these lengths is consistent with the lengths of the particular polynucleotide, of a polynucleotide described in Table III or Table IV, or a sequence complementary thereto, wherein said
- 15 polynucleotide described in Table III or Table IV is selected individually or in any combination from the polynucleotides described in Table III or Table IV. The present invention also encompasses isolated, purified, or recombinant nucleic acids which consist of or consist essentially of a polynucleotide described in Table III or Table IV, or a sequence complementary thereto, wherein said polynucleotide is selected individually or in any combination from the polynucleotides described in Table III or Table IV.

20

**Table III**

<b>SEQ ID NO</b>	<b>Positions of preferred fragments</b>
24	1-251
25	1-83
28	227-276
29	1-27
30	130-242, 283-315, 365-461
32	314-399
33	89-321
34	1-38
35	1-52, 171-222
36	1-30, 408-441
37	1-138
39	115-140
40	1-97
41	1-112
42	1-177
46	1-38
48	376-400
51	400-466
54	1-259
55	189-320

56	265-457
58	246-469
59	81-123, 418-444
60	1-348
61	78-123, 418-457
62	386-439
63	1-214
64	109-297
65	1-370
66	92-428
68	1-180
69	165-259
70	1-178
71	1-27
72	1-179
73	1-65, 107-192
75	1-314
77	263-388
78	1-64
79	1-149
80	101-142, 302-380
82	1-192
83	1-398
85	1-290
86	1-118, 149-336
87	1-262
88	1-149
89	1-315
90	1-74
91	1-335, 364-423
92	1-316
93	338-508
94	179-321
95	219-402
96	26-315
97	348-460
98	1-230
99	391-467
101	214-336
102	1-289
103	1-383
104	1-211
105	1-36
106	1-126
107	1-49
108	294-336
109	1-128
111	1-154
112	407-441
113	1-80, 139-184
114	10-79
116	1-292
117	1-304

119	1-288
120	2-348
121	1-122
123	188-353
124	1-249
125	295-375
128	1-244
129	1-232
130	196-312
131	178-276
132	37-174
133	1-344
134	1-244
135	1-217
136	82-428
137	1-29, 103-155, 274-434
138	1-395
139	1-268
140	1-170
141	1-396
142	1-73, 227-357
143	1-159
144	1-433
145	61-116
146	1-71, 179-205
147	177-300
149	1-146
151	1-166
152	1-382
153	1-208
154	121-251
155	1-147
157	1-115
158	1-175
159	1-44, 80-230
160	1-346
161	1-277
162	1-235
163	1-34
164	1-195
165	19-78, 175-217
166	1-209
167	1-65
168	128-218
169	49-245
170	179-280
171	1-103
172	1-218
173	1-380
174	1-139
175	1-122
176	1-300
177	1-466

179	1-86
180	1-245
181	1-241
182	1-263
183	1-170
184	58-106, 399-443
185	1-427
186	1-365
187	1-260
188	1-172
189	1-150
190	161-271, 301-339
191	1-91
192	1-264
193	1-246
194	1-150
195	1-209
196	1-363
197	1-155
198	1-135
200	1-125
201	1-210
202	1-338
203	1-188
204	228-347
205	1-440
206	56-221
208	1-422
209	169-195
210	1-363
211	1-368
212	1-448
213	1-134
214	1-193
215	1-214
216	1-134
218	1-189
219	1-248
220	1-115
221	1-113
222	1-370
224	1-251
225	1-198
226	45-141
227	1-206
228	1-480
229	1-144
230	1-42, 281-351, 432-457
231	1-112
233	1-301
234	1-109
235	1-393
236	1-222

237	1-154
238	1-439
239	112-137
240	1-194
241	1-44
242	1-242
244	1-324
245	1-38, 217-280
246	1-60
247	77-359
248	1-236
249	1-342
250	80-382
251	1-303
252	62-259
253	1-165
254	1-328
255	1-320
256	1-305
257	1-181
258	116-174
259	1-265
260	1-272
261	1-62
263	1-371
266	1-274
267	1-342
268	364-427
269	31-143
270	1-79
271	1-121
272	229-292
273	1-158
274	1-113
275	1-254
276	1-333
277	1-130
278	1-184
279	1-265
280	1-188
281	1-177
282	1-336
283	1-294
284	1-171
285	1-297
288	1-42
290	1-170
292	20-155
294	1-334
295	1-375
296	1-226
297	1-232
299	40-139

300	1-285
301	1-242
302	1-136
303	1-175
304	1-493
305	1-214
306	89-458
307	1-328
308	1-380
309	1-236
310	1-357
311	1-470
312	1-187
313	1-159
315	1-162
316	1-404
317	1-450
318	1-395
319	1-257
320	56-325
321	1-201
322	1-159
323	1-420
324	1-210
325	1-192
326	88-181
327	1-185
328	128-210
330	1-223
331	1-362
332	1-89
334	1-188
335	1-115
336	1-300
337	1-307
338	1-123
339	1-297
340	1-34
341	1-44
342	1-37
343	141-169
344	1-112
345	1-235, 266-349
346	1-191
347	1-229
348	1-210
350	139-266
351	1-307
352	1-170
353	1-293
354	30-161, 192-331
355	1-93
356	1-178



357	1-107
358	1-29, 168-209
359	1-298
360	1-193
362	1-360
363	1-45, 100-212
364	39-170, 202-242
365	1-248
366	1-351
367	1-208
368	228-446
369	1-62
370	1-132
371	1-127
372	1-196
373	1-148
374	1-126
375	1-112
376	1-146
378	1-143
379	1-261
380	202-228
382	1-151
383	1-45
384	1-190, 250-456
385	1-55, 141-181
386	1-281
387	1-111
388	1-374
389	1-192
390	1-371
392	1-303
394	1-126
395	1-329
396	1-99
397	1-316
398	1-251
399	1-120
401	1-206
402	1-330
403	1-311
405	1-153
406	1-206
407	1-479
408	1-289
410	229-321
413	1-158
415	95-229
416	1-265
417	1-228
418	1-225
419	207-293
420	1-194

421	1-90
422	1-161
423	1-420
424	1-432
425	1-276, 309-419
426	1-232
427	1-81
428	1-96
429	1-165
431	1-58, 186-237, 327-354
433	1-65
434	1-83
435	1-386
436	405-447
438	1-106
439	45-105, 168-255, 284-447
441	1-409
442	1-320
443	1-256
444	1-284
445	1-240
446	1-149
447	1-360
448	1-123
449	1-94
450	1-302
452	1-349
453	1-270
454	1-492
455	17-105
456	1-102
457	1-108
458	1-285
459	1-311
460	1-191
461	312-420
462	1-257
463	1-117
464	1-142
466	1-235
467	1-29
468	1-41
469	1-438
470	1-131
471	1-211
472	1-150
473	1-352
474	1-141
476	1-232
478	1-201
479	1-151
480	1-104
481	7-429

482	1-385
486	1-226
488	1-296
489	1-72, 323-377
491	1-348
492	33-126
493	1-300
494	1-295
495	1-244
496	1-215
497	1-255
499	1-174, 384-474
500	1-50, 102-241
501	153-430
502	91-132
503	1-64
504	21-63, 356-420
505	37-68, 187-234
506	1-315
507	101-208
510	1-402
511	1-343
512	1-140, 170-246, 276-420
513	1-324
514	1-303
515	13-340
516	1-263, 293-360
518	1-245
519	111-275
520	62-182
521	1-218
523	1-502
524	1-118
525	1-276
526	223-366
527	1-428
528	297-342
529	1-244
530	1-88, 375-434
531	1-406
533	1-149
534	1-145
535	1-116
536	1-207
537	1-394
538	1-415
539	1-160
540	1-327
541	1-38, 73-396
542	1-247
543	1-221
544	1-375
545	1-376

546	1-109
547	1-160, 223-306
548	1-148
551	1-231
552	1-229
553	1-232
554	1-141
555	1-376
556	1-279
557	1-340
558	1-51
559	1-354
562	1-188
563	1-229
564	184-352
566	308-341
567	1-218
568	1-79
569	1-142
570	1-207
571	1-373
572	1-195
573	1-352
574	1-121
575	1-222
576	151-288
577	1-264
578	1-205
580	1-171, 273-328
581	1-356
582	1-239
583	1-144
584	1-282
585	1-338
586	1-436
588	1-380
589	1-60
590	1-178
592	1-66
593	1-215
594	1-161
596	1-407
597	31-83
598	1-417
599	1-329
600	1-311
601	1-61, 99-214
602	1-154, 197-463
603	135-269
604	1-351
605	1-195
608	1-357
609	1-201

612	1-176
613	1-342
615	1-272
616	1-114
617	1-46
618	1-208
619	1-257
620	1-28
621	1-26
622	1-221
623	1-432
624	1-233
625	1-26
627	1-43
628	1-318
629	1-170
630	1-196
631	248-339
632	1-433
633	1-154
634	1-41
635	1-137
636	1-172
637	1-253
638	1-185
639	1-206
641	334-483
642	1-309
643	1-75, 162-213
644	107-211
645	1-98
646	1-347
647	1-49, 81-143
648	1-232
649	74-133
650	1-37
651	1-276
652	1-170
653	1-178
654	1-121
656	1-197
657	1-246
659	1-197
660	116-172
661	1-411
662	1-146
663	1-65
664	1-182
665	1-320
666	1-273
667	1-149
668	1-122
670	1-160

671	1-137
673	1-263
674	1-263
675	1-107
677	1-441
678	134-191
679	1-235
680	1-26
682	1-58, 269-328
683	1-447
684	1-217
685	1-132
686	1-60
688	1-107
689	132-221, 327-377
690	1-388
691	1-141, 171-408
692	1-322
693	1-153
695	1-455
698	1-58, 117-174
699	240-300
700	1-159
701	1-69
702	1-175
703	1-298
704	1-136
705	1-168
706	1-419
707	1-382
708	8-245, 296-384
709	1-149
710	1-167
711	1-35
712	1-80, 116-156, 206-241
713	33-376
714	1-304
715	1-242
717	1-145
718	1-350
720	1-257
721	1-360
722	1-191
724	1-139
726	1-207
727	99-164
728	1-321
730	156-372
731	1-109, 256-290
735	25-192
737	1-160
738	1-227
739	441-514

742	217-280
743	10-275
747	1-179
749	2-31, 139-168
750	349-410
752	1-119
753	1-121
754	1-28
760	25-175
761	1-212
763	8-75
766	1-59, 102-248, 295-320
769	53-85
771	1-370
774	1-347
776	1-200
778	39-342
779	4-28
780	1-49, 407-472
781	116-426
782	1-59
783	1-53, 219-453
784	29-53, 219-263, 426-494
785	99-347, 386-461
786	2-28
788	1-279
789	1-58
790	226-268
792	129-218
794	265-431
796	5-86
797	1-34
799	1-344
802	46-477
806	64-384
807	135-301
808	2-314
810	6-39
1600	1-25
1601	1-225
1602	23-139
1603	1-294
1606	15-44
1607	1-361
1611	85-228
1612	1-221
1613	138-281
1614	65-171
1615	2-142
1616	1-46
1617	1-95
1620	1-187
1621	1-136

1622	32-280, 311-400
------	-----------------

5

Table IV

SEQ ID NO	Positions of Preferred Fragments
35	1-52
41	1-115
45	1-47
46	1-33
66	400-428
82	83-149
93	399-508
105	1-36
114	1-79
120	1-386
121	1-190
124	1-249
125	295-328
139	1-81, 125-268
159	1-139, 180-230
165	1-78
170	179-205, 248-280
194	1-150
213	1-158
247	1-104, 155-183, 280-359
269	31-143
350	139-386
368	228-446
385	1-72, 143-193
415	95-229
435	1-386
436	446-472
441	1-361
454	1-349
455	1-105
459	35-161, 200-311
460	1-26, 56-140
481	1-429
489	1-84
496	1-44, 84-215
501	153-430
502	1-91
504	1-63



505	1-68
514	1-303
515	237-351
519	1-145
526	231-366
530	1-88
535	1-55
570	76-207
576	168-218, 261-288
588	1-331
597	1-83
627	1-43
634	1-41
641	1-55, 334-483
672	1-34
687	1-129
708	1-245, 296-384
710	1-26, 104-167
722	1-191
730	1-465
731	1-43
735	1-91
737	1-160
738	1-186
739	1-48
742	1-62, 99-248
743	1-315, 412-459
744	1-31
747	1-63
749	1-32
750	1-38
752	1-139
753	1-193
754	1-28
759	1-38
760	1-115
763	1-62
765	1-126
769	1-85
770	1-40
771	1-148
774	1-134
775	265-531
776	71-203
777	333-469
778	144-468
779	1-28
780	1-49
781	1-102
782	1-59
783	1-53
784	1-220, 262-390
785	1-339, 408-461

786	1-28
789	1-58
791	1-126
792	1-31, 129-220
793	1-31
794	355-431
795	1-33
797	1-31
798	1-31
799	1-401
801	1-117
802	1-92
806	64-384
807	1-331
808	1-351
810	1-39
1600	1-25
1603	1-341
1606	1-31
1607	1-361
1608	164-305
1611	85-228
1612	1-221
1613	112-360
1614	1-171
1615	94-193
1617	1-155
1620	1-246

### III. Evaluation of Spatial and Temporal Expression of mRNAs Corresponding to the 5'ESTs, Consensus Contigated 5'ESTs, or EST-related nucleic acids

5

#### EXAMPLE 17

##### Expression Patterns of mRNAs From Which the 5'ESTs were obtained

Each of the SEQ ID NOs. 24-811 and 1600-1622 was also categorized based on the tissue from which its corresponding mRNA was obtained, as follows.

10 Table V shows the spatial distribution of each nucleic acid sequence of the invention (SEQ ID NOs. 24-811 and 1600-1622) referred to by its sequence identification number in the first column. In the second column entitled tissue distribution, the spatial distribution is represented by the number of individual 5'ESTs used to assemble the consensus contigated 5'ESTs for a given tissue. Each type of tissue listed in Table V is encoded by a letter. The correspondence between the letter code and the tissue  
15 type is given in Table VI.

Table V

SEQ ID NO	Tissue Distribution
24	AA:1
25	S:1
26	P:1
27	W:1
28	P:1
29	S:1
30	P:1
31	P:1
32	P:1
33	P:1
34	AB:1
35	G:3; P:1; S:1; W:3; AA:4
36	P:1
37	S:1
38	Q:1
39	P:1
40	AB:1
41	B:1; C:3; F:1; G:1; H:4; S:2; T:8; W:1; Z:1; AA:3; AC:1; AD:3
42	A:1
43	N:2
44	P:1
45	C:2; K:1; O:1; S:5
46	K:1; S:2; AA:1
47	AA:1
48	C:1; O:1; P:8
49	P:1
50	P:1
51	P:1
52	S:1
53	AA:1
54	T:1
55	P:1
56	P:1
57	P:1
58	P:1
59	P:7; T:2; Z:1
60	R:1
61	C:1
62	P:1
63	F:1
64	AA:1
65	F:1

66	P:4; T:2; Z:1
67	S:1
68	AA:1
69	P:1
70	P:1
71	S:1
72	W:1
73	G:1
74	P:1
75	N:1
76	P:1
77	S:1
78	U:1
79	B:1
80	P:1
81	AC:1
82	K:1; O:1
83	G:1
84	C:1; K:2; P:29; S:2; T:1; X:2; Y:1; AA:2
85	K:1
86	C:1
87	F:1
88	AB:1
89	H:1
90	M:1
91	B:1
92	K:1
93	AC:2
94	P:1
95	M:1
96	Z:2
97	K:1; P:11; S:1; X:1; AA:1
98	W:1
99	X:1
100	P:1
101	AB:1
102	F:1
103	AA:1
104	K:1
105	B:4; C:6; E:2; H:3; O:2; Q:1; S:3; AC:2
106	T:1
107	O:1
108	P:1
109	G:1
110	AA:1
111	T:1
112	P:1
113	F:1

114	B:3; C:4; K:5; S:4; Y:1
115	U:1
116	W:1
117	T:1
118	T:2
119	T:1
120	H:3
121	AA:3
122	K:1
123	H:2
124	AA:2
125	B:1; G:1; J:3; T:13; Y:5; AA:5; AD:2
126	H:1; P:1
127	K:1
128	F:1
129	G:1
130	P:1
131	B:1
132	AA:1
133	W:1
134	P:1
135	K:1
136	B:1; C:1
137	B:1
138	H:1
139	AC:2
140	T:1
141	B:1
142	H:1
143	T:1
144	H:1
145	B:1
146	R:1
147	P:1
148	C:1; H:2; O:1; S:2; T:1; AC:2
149	H:1
150	AA:1
151	W:1
152	S:1
153	F:1
154	M:1
155	B:1
156	R:1
157	W:1
158	T:1
159	C:1; AA:1
160	F:1
161	H:1

162	D:1
163	AA:1
164	AA:1
165	W:3
166	AA:1
167	W:1
168	F:1
169	B:1
170	G:2
171	E:1
172	B:1
173	F:1
174	B:1
175	W:1
176	K:1
177	AA:1
178	S:1
179	K:1
180	AA:1
181	W:1
182	K:1
183	T:1
184	P:1
185	B:1
186	W:1
187	R:1
188	T:1
189	T:1
190	W:1
191	A:1
192	F:1
193	B:1
194	G:3
195	W:1
196	O:1
197	T:1
198	O:1
199	B:1
200	AA:1
201	G:1
202	B:1
203	G:1
204	P:1
205	AA:1
206	Y:1
207	Y:1
208	AA:1
209	G:1

210	H:1
211	C:1
212	H:1
213	W:2
214	Y:1
215	AB:1
216	K:1
217	M:1
218	AD:1
219	A:1
220	AA:1
221	G:1
222	G:1
223	G:1; H:2; S:2; X:1
224	G:1
225	G:1
226	B:1
227	P:1
228	O:1
229	G:1
230	T:1
231	T:1
232	K:1
233	S:1
234	O:1
235	F:1
236	T:1
237	B:1
238	W:1
239	G:1
240	R:1
241	A:1
242	W:1
243	P:1
244	H:1
245	D:1
246	C:1
247	B:2
248	P:1
249	F:1
250	AB:1
251	W:1
252	H:1
253	B:1
254	S:1
255	T:1
256	W:1
257	T:1

258	AA:2
259	P:1
260	W:1
261	H:1
262	K:1
263	K:1
264	C:1; E:1; F:1; I:4; L:1; N:22; O:1; P:1; S:1; T:9; AA:1
265	A:1
266	T:1
267	K:1
268	H:1
269	T:2
270	T:1
271	T:1
272	B:1
273	Y:1
274	T:1
275	G:1
276	AA:1
277	T:1
278	AB:1
279	T:1
280	W:1
281	F:1
282	K:1
283	H:1
284	O:1
285	W:1
286	B:21; C:7; H:5; K:5; O:8; S:16; W:1; Y:3; Z:4; AA:2; AC:1
287	K:2; P:12; W:1; AC:2
288	S:1
289	K:2; P:8; W:1; AC:2
290	S:1
291	H:1
292	B:11; C:2; E:1; H:7; K:1; N:3; S:1; T:8; W:1; AA:28; AC:1
293	B:6; C:3; G:1; H:6; K:4; N:4; O:3; Q:2; S:5; T:1; U:1; V:2; Y:3; AA:1
294	B:1
295	H:1
296	AA:1
297	T:1
298	T:1
299	T:1
300	H:1; S:1
301	H:1
302	W:1
303	W:1
304	H:1
305	G:1



306	K:1
307	H:1
308	A:1
309	H:1
310	H:1
311	Y:1
312	G:1
313	H:1
314	K:1
315	Y:1
316	P:1
317	H:1
318	AA:1
319	H:1
320	O:1
321	Y:1
322	B:1
323	P:1
324	P:1
325	K:1
326	H:1
327	H:1
328	Q:1
329	S:1
330	B:1
331	T:1
332	T:1
333	B:1
334	T:1
335	W:1
336	P:1
337	A:1
338	AA:1
339	AA:1
340	G:1
341	C:1
342	K:1
343	S:1
344	G:1
345	B:1
346	Y:1
347	G:1
348	F:1
349	AA:5
350	B:15; C:1; G:1; H:1; O:1; Q:2; S:1; X:1; Y:1
351	F:1
352	R:1
353	O:1

354	H:1
355	W:1
356	F:1
357	T:1
358	S:1
359	X:1
360	T:1
361	K:1
362	K:1
363	G:1
364	K:1
365	G:1
366	AA:1
367	F:1
368	C:2; H:2; X:1
369	E:1
370	T:1
371	H:1
372	G:1
373	AA:1
374	G:1
375	F:1
376	F:1
377	R:1
378	AA:1
379	AA:1
380	C:1
381	H:1
382	T:1
383	W:1
384	S:1
385	AA:2
386	D:1
387	O:1
388	W:1
389	F:1
390	W:1
391	K:1
392	W:1
393	K:1
394	T:1
395	H:1
396	T:1
397	T:1
398	G:1
399	C:1
400	K:1
401	B:1

402	H:1
403	B:1
404	B:1
405	H:1
406	AB:1
407	O:1
408	P:1
409	X:1
410	H:1
411	B:9; C:3; K:3; L:2; O:1; S:2; X:1; AA:1
412	G:1; S:2; V:2; W:1; Y:1; Z:1
413	W:1
414	G:1
415	B:3; C:3; F:1; G:2; H:4; J:1; K:1; O:1; P:3; S:1; V:1
416	I:1
417	F:1
418	F:1
419	F:1
420	AA:1
421	F:1
422	T:1
423	P:1
424	B:1
425	Y:1
426	W:1
427	AA:1
428	W:1
429	H:1
430	Y:1
431	J:1
432	AA:1
433	G:1
434	AA:1
435	B:3; H:1
436	B:9; G:4; H:8; K:2; O:2; W:1; Z:2; AA:2; AD:3
437	H:1; T:1
438	T:1
439	R:1
440	M:1
441	H:2
442	W:1
443	B:1
444	W:1
445	AB:1
446	F:1
447	AD:1
448	AB:1
449	N:1

450	T:1
451	W:1
452	O:1
453	AA:1
454	D:28
455	W:1
456	T:1
457	G:1
458	W:1
459	Y:4
460	B:3
461	P:2
462	K:1
463	T:1
464	H:1
465	G:1
466	AC:1
467	R:1
468	S:1
469	B:1
470	S:1
471	T:1
472	AA:1
473	W:1
474	T:1
475	S:1
476	T:1
477	AA:1
478	G:1
479	W:1
480	B:1
481	O:2
482	K:1
483	P:1
484	W:1
485	P:1
486	B:1
487	Y:1
488	H:1
489	P:1; Q:1; S:3
490	C:1
491	S:1
492	H:1
493	B:1
494	H:1
495	G:1
496	N:2
497	B:1

498	G:1
499	P:1
500	G:1
501	C:1; K:1; Q:1
502	B:4
503	R:1
504	B:5; H:2; W:2
505	G:2; H:1
506	W:1
507	B:1
508	W:1
509	AB:1
510	H:1
511	N:1
512	J:1
513	AA:1
514	T:2
515	AA:5
516	F:1
517	C:1; O:1
518	W:1
519	T:4
520	B:1
521	H:1
522	H:2; T:3
523	H:1
524	AA:1
525	W:1
526	C:2; E:1; J:1; R:3; S:4; AA:1
527	H:1
528	S:1
529	P:1
530	B:1; H:1
531	O:1
532	Y:1
533	H:1
534	T:1
535	T:2
536	B:1
537	AD:1
538	AA:1
539	T:1
540	F:1
541	AD:1
542	W:1
543	W:1
544	F:1
545	T:1

546	F:1
547	K:1
548	Y:1
549	S:1
550	B:1
551	B:1
552	B:1
553	H:1
554	P:1
555	G:1
556	H:1
557	K:1
558	B:1
559	R:1
560	AB:1
561	C:1; S:1; V:1
562	AA:1
563	K:1
564	P:1
565	K:1
566	G:1
567	W:1
568	E:1; W:2
569	W:1
570	B:2
571	O:1
572	T:1
573	B:1
574	T:1
575	B:1
576	B:3
577	B:1
578	X:1
579	H:1
580	AA:1
581	AA:1
582	AA:1
583	AA:1
584	AA:1
585	D:1
586	H:1
587	H:1
588	AA:3
589	K:1
590	W:1
591	K:1
592	W:1
593	B:1

594	V:1
595	R:1
596	P:1
597	G:1; X:2; Z:1
598	X:1
599	F:1
600	F:1
601	Y:1
602	F:1
603	W:1
604	H:1
605	G:1
606	C:2; H:1; S:3; W:2; AD:3
607	W:1
608	C:1
609	F:1
610	K:1
611	M:1
612	AD:1
613	H:1
614	T:1
615	H:1
616	F:1
617	T:1
618	G:1
619	G:1
620	B:1
621	W:1
622	W:1
623	T:1
624	AA:1
625	G:1
626	M:1
627	C:2; T:2; W:1; Y:1
628	T:1
629	J:1
630	T:1
631	P:1
632	H:1
633	H:1
634	C:1; S:1; T:1; AD:1
635	J:1
636	G:1
637	W:1
638	AA:1
639	W:1
640	B:6; C:3; G:1; H:2; K:6; O:4; Q:1; R:2; S:1; T:3; Y:3; Z:2; AA:2; AC:2; AD:3

641	B:21; C:2; G:5; W:4; Y:1
642	AA:1
643	P:1
644	AA:1
645	T:1
646	K:1
647	F:1
648	F:1
649	F:1
650	T:1
651	W:1
652	T:1
653	T:1
654	P:1
655	B:1; H:2; N:1; T:3; Y:1
656	B:1
657	T:1
658	R:1
659	K:1
660	W:1
661	AA:1
662	Y:1
663	W:1
664	G:1
665	S:1
666	Y:1
667	F:1
668	T:1
669	B:1
670	F:1
671	T:1
672	A:2; B:6; C:1; G:1; H:3; J:1; L:1; P:2; Q:1; S:4; T:1; V:3; W:2; Y:1; AA:3; AD:2
673	T:1
674	G:1
675	F:1
676	M:1
677	G:1
678	Y:1
679	D:1
680	P:1
681	D:1
682	AA:1
683	G:1
684	K:1
685	G:1
686	P:1
687	B:3; C:2; D:2; E:2; J:4; V:2; AC:6



688	AA:1
689	S:1
690	AA:1
691	H:1
692	AA:1
693	S:1
694	AB:1
695	T:1
696	H:1
697	B:4; E:1; F:1; P:1; T:2; Z:2
698	O:1
699	W:1
700	S:1
701	O:1
702	B:1
703	AB:1
704	H:1
705	B:1
706	H:1
707	G:1
708	F:1; H:1; K:1; W:2; AA:1
709	H:1
710	T:2
711	C:1
712	G:1
713	Y:1
714	C:1
715	Y:1
716	Z:1
717	P:1
718	G:1
719	S:1
720	K:1
721	M:1
722	T:2
723	O:1; P:2; S:2
724	T:1
725	T:1
726	N:1
727	T:1
728	T:1
729	C:2; H:2; K:2; V:1; AC:1
730	B:7; H:2; Y:1
731	B:5; W:3
732	B:1; C:2; G:2; S:2; AA:9
733	B:6; C:2; G:1; H:10; O:2; P:6; Q:1; S:2; W:4; AC:2
734	B:6; O:1; V:1
735	C:1; O:2

736	B:1; H:2; N:1; T:3; Y:1
737	T:2
738	T:2
739	B:3; C:8; D:1; E:6; G:3; H:11; I:1; J:1; N:1; O:3; P:12; Q:3; S:2; T:2; W:1; AC:1; AD:8
740	H:2; Y:1
741	C:2; H:1
742	B:12; C:1; G:1; H:4; K:2; O:2; S:4; T:2; Y:2
743	AA:4
744	B:1; G:1; H:6; T:1; W:1
745	C:7; E:1; G:3; H:2; P:2; S:2; T:1; W:1; AD:2
746	G:2; S:1
747	T:2
748	S:3
749	H:1; O:2; S:2
750	Y:1; AD:1
751	B:8; G:2; H:2; I:1; Q:2; S:2; T:1; W:2
752	T:3
753	P:4
754	B:1; H:2
755	B:7; C:1; G:6; H:2; K:1; U:2; V:1; Z:1
756	C:1; H:1; J:2; O:2; S:1; T:2; W:1; AA:1
757	B:1; C:1; K:3; S:1; V:1; Y:1
758	E:1; H:2; K:1; P:1; Q:1; AD:5
759	B:6; C:1; Y:1
760	B:4
761	W:2
762	B:3; C:7; H:9; N:1; S:1; T:1; Y:1; AA:1
763	N:1; S:1; AA:5
764	H:3
765	B:3; G:1; W:1
766	H:2
767	C:1; AA:3
768	B:2; C:6; H:9; N:1; S:1; T:1; Y:1; AA:1
769	A:1; B:4; C:4; F:4; G:6; H:10; K:2; O:8; P:2; R:1; S:8; T:2; W:3; AA:2; AC:1
770	A:2; P:16; X:1
771	AA:3
772	O:4
773	B:1; C:1; W:1
774	P:2; X:4
775	B:18; C:6; H:5; K:3; O:7; S:10; W:1; Y:3; Z:2; AA:2; AC:1
776	H:7
777	B:26; C:8; H:5; K:4; O:10; S:17; W:1; Y:4; Z:4; AA:4; AC:2
778	B:6
779	B:3; C:1; G:1; H:2; K:1; Q:1; S:8; W:2; Y:9; AA:4
780	B:3; C:1; F:1; P:1; W:1; AC:1
781	I:2; N:1; P:1; R:3; AA:1
782	B:2

783	H:1; P:2; S:3; AD:1
784	H:1; P:1; S:4; AD:1
785	T:2
786	D:1; AC:9
787	H:1; L:1; S:1
788	B:6; S:4
789	S:1; T:1
790	B:1; C:2; H:5; W:1; AD:1
791	B:3; C:2; D:3; E:2; J:4; V:3; AC:5
792	B:3; D:1; K:2; S:2; Y:1
793	B:2; G:2; AA:1
794	B:25; C:4; D:1; E:1; F:3; G:6; J:1; K:6; N:1; O:1; P:2; R:1; S:3; T:2; W:2; X:1; Y:1; Z:1; AA:1; AC:2; AD:1
795	B:4; C:1; E:2; H:4; J:1; L:1; O:4; S:1; V:1; Y:3; Z:1
796	H:5
797	B:2; E:1; N:2
798	B:1; G:1; H:6; T:1; W:1
799	H:2
800	H:2; I:2; AA:1
801	A:2; B:4; C:14; D:1; H:2; K:1; N:2; S:4; T:1; W:2; AA:20
802	AA:17
803	B:2; G:3; H:3; S:1; U:1; AC:1; AD:2
804	C:1; S:2; T:2; X:2; AA:1; AC:1
805	B:5; C:6; D:5; H:17; J:2; K:4; N:1; O:6; P:2; S:5; T:5; W:1; X:1; Z:2; AA:13; AC:3
806	B:2; C:3; D:3; H:6; J:2; K:1; N:1; O:3; P:1; S:2; T:4; W:1; X:1; Z:1; AA:5; AC:1
807	H:1; AC:4
808	R:13
809	B:3; W:4
810	B:16; S:1; Y:14
811	B:8; C:5; G:1; H:1; K:5; O:2; Q:2; R:2; S:2; T:3; Y:4; Z:2; AA:1; AC:1; AD:2
1600	T:4
1601	AA:3
1602	C:3; H:1
1603	H:2; AC:2
1604	B:7; C:1; E:1; H:1; P:2; R:3; S:2; T:2; Z:3; AA:2
1605	C:4; H:3; O:1
1606	A:3; B:13; C:14; D:2; E:10; F:3; G:19; H:32; K:11; O:5; P:2; R:3; S:16; T:4; W:2; Y:10; Z:8; AA:1; AC:3
1607	T:3
1608	B:3; P:2
1609	R:4
1610	B:4
1611	B:3; T:1
1612	T:2
1613	V:5
1614	D:3

1615	AA:10
1616	B:4
1617	T:2
1618	K:2; S:8; AA:1
1619	B:2
1620	W:2
1621	H:1; AB:1
1622	H:2

Table VI

Tissue code	Tissue type
A	Bone Marrow
B	Brain
c	Cancerous prostate
D	Cerebellum
E	Colon
F	Dystrophic muscle
G	Fetal brain
H	Fetal kidney
I	Fetal liver
J	Heart
K	Hypertrophic prostate
L	Kidney
M	Large intestine
N	Liver
O	Lung
P	Lymph ganglia
Q	Lymphocytes
R	Muscle
S	Prostate
T	Ovary
U	Pancreas
V	Placenta
W	Spinal cord
X	Spleen
Y	Substantia nigra
Z	Surrenals
AA	Testis
AB	Thyroid
AC	Umbilical cord
AD	Uterus

- 5 In addition to categorizing the 5' ESTs and consensus contigated 5' ESTs with respect to their tissue of origin, the spatial and temporal expression patterns of the mRNAs corresponding to the 5' ESTs and consensus contigated 5' ESTs, as well as their expression levels, may be determined as described in Example 18 below.

Characterization of the spatial and temporal expression patterns and expression levels of these mRNAs is useful for constructing expression vectors capable of producing a desired level of gene product in a desired spatial or temporal manner, as will be discussed in more detail below.

Furthermore, 5' ESTs and consensus contigated 5' ESTs whose corresponding mRNAs are associated with disease states may also be identified. For example, a particular disease may result from the lack of expression, over expression, or under expression of a mRNA corresponding to a 5' EST or consensus contigated 5' EST. By comparing mRNA expression patterns and quantities in samples taken from healthy individuals with those from individuals suffering from a particular disease, 5' ESTs or consensus contigated 5' ESTs responsible for the disease may be identified.

It will be appreciated that the results of the above characterization procedures for 5' ESTs and consensus contigated 5' ESTs also apply to extended cDNAs (obtainable as described below) which contain sequences adjacent to the 5' ESTs and consensus contigated 5' ESTs. It will also be appreciated that if desired, characterization may be delayed until extended cDNAs have been obtained rather than characterizing the 5' ESTs or consensus contigated 5' ESTs themselves.

### EXAMPLE 18

#### Evaluation of Expression Levels and Patterns of mRNAs

##### Corresponding to EST-Related Nucleic Acids

Expression levels and patterns of mRNAs corresponding to EST-related nucleic acids may be analyzed by solution hybridization with long probes as described in International Patent Application No. WO 97/05277. Briefly, an EST-related nucleic acid, fragment of an EST-related nucleic acid, positional segment of an EST-related nucleic acid, or fragment of a positional segment of an EST-related nucleic acid corresponding to the gene encoding the mRNA to be characterized is inserted at a cloning site immediately downstream of a bacteriophage (T3, T7 or SP6) RNA polymerase promoter to produce antisense RNA. Preferably, the EST-related nucleic acid, fragment of an EST-related nucleic acid, positional segment of an EST-related nucleic acid, or fragment of a positional segment of an EST-related nucleic acid is 100 or more nucleotides in length. The plasmid is linearized and transcribed in the presence of ribonucleotides comprising modified ribonucleotides (*i.e.* biotin-UTP and DIG-UTP). An excess of this doubly labeled RNA is hybridized in solution with mRNA isolated from cells or tissues of interest. The hybridizations are performed under standard stringent conditions (40-50°C for 16 hours in an 80% formamide, 0.4 M NaCl buffer, pH 7-8). The unhybridized probe is removed by digestion with ribonucleases specific for single-stranded RNA (*i.e.* RNases CL3, T1, Phy M, U2 or A). The presence of the biotin-UTP modification enables capture of the hybrid on a microtitration plate coated with streptavidin. The presence of the DIG modification enables the hybrid to be detected and quantified by ELISA using an anti-DIG antibody coupled to alkaline phosphatase.

The EST-related nucleic acid, fragment of an EST-related nucleic acid, positional segment of an EST-related nucleic acid, or fragment of a positional segment of an EST-related nucleic acid may also be

tagged with nucleotide sequences for the serial analysis of gene expression (SAGE) as disclosed in UK Patent Application No. 2 305 241 A. In this method, cDNAs are prepared from a cell, tissue, organism or other source of nucleic acid for which gene expression patterns must be determined. The resulting cDNAs are separated into two pools. The cDNAs in each pool are cleaved with a first restriction  
5 endonuclease, called an anchoring enzyme, having a recognition site which is likely to be present at least once in most cDNAs. The fragments which contain the 5' or 3' most region of the cleaved cDNA are isolated by binding to a capture medium such as streptavidin coated beads. A first oligonucleotide linker having a first sequence for hybridization of an amplification primer and an internal restriction site for a so called tagging endonuclease is ligated to the digested cDNAs in the first pool. Digestion with the  
10 second endonuclease produces short tag fragments from the cDNAs.

A second oligonucleotide having a second sequence for hybridization of an amplification primer and an internal restriction site is ligated to the digested cDNAs in the second pool. The cDNA fragments in the second pool are also digested with the tagging endonuclease to generate short tag fragments derived from the cDNAs in the second pool. The tags resulting from digestion of the first and second  
15 pools with the anchoring enzyme and the tagging endonuclease are ligated to one another to produce so called ditags. In some embodiments, the ditags are concatamerized to produce ligation products containing from 2 to 200 ditags. The tag sequences are then determined and compared to the sequences of the EST-related nucleic acid, fragment of an EST-related nucleic acid, positional segment of an EST-related nucleic acid, or fragment of a positional segment of an EST-related nucleic acid to determine  
20 which 5' ESTs, consensus contigated 5' ESTs, or extended cDNAs are expressed in the cell, tissue, organism, or other source of nucleic acids from which the tags were derived. In this way, the expression pattern of the 5' ESTs, consensus contigated 5' ESTs, or extended cDNAs in the cell, tissue, organism, or other source of nucleic acids is obtained.

Quantitative analysis of gene expression may also be performed using arrays. As used herein,  
25 the term array means a one dimensional, two dimensional, or multidimensional arrangement of EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments EST-related nucleic acids, or fragments of positional segments of EST-related nucleic acids. Preferably, the EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments EST-related nucleic acids, or fragments of positional segments of EST-related nucleic acids are at least 10, 12, 15, 18, 20, 23, 25, 28,  
30 30, 35, 40, or 50 nucleotides in length. More preferably, the EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments EST-related nucleic acids, or fragments of positional segments of EST-related nucleic acids are at least 100 nucleotide long. More preferably, the fragments are more than 100 nucleotides in length. In some embodiments, the EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments EST-related nucleic acids, or fragments of  
35 positional segments of EST-related nucleic acids may be more than 500 nucleotides long.

For example, quantitative analysis of gene expression may be performed with EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments EST-related nucleic acids, or

fragments of positional segments of EST-related nucleic acids in a complementary DNA microarray as described by Schena *et al.* (*Science* 270:467-470, 1995; *Proc. Natl. Acad. Sci. U.S.A.* 93:10614-10619, 1996). EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments EST-related nucleic acids, or fragments of positional segments of EST-related nucleic acids are amplified by  
5 PCR and arrayed from 96-well microtiter plates onto silylated microscope slides using high-speed robotics. Printed arrays are incubated in a humid chamber to allow rehydration of the array elements and rinsed, once in 0.2% SDS for 1 min, twice in water for 1 min and once for 5 min in sodium borohydride solution. The arrays are submerged in water for 2 min at 95°C, transferred into 0.2% SDS for 1 min, rinsed twice with water, air dried and stored in the dark at 25°C.

10 Cell or tissue mRNA is isolated or commercially obtained and probes are prepared by a single round of reverse transcription. Probes are hybridized to 1 cm<sup>2</sup> microarrays under a 14 x 14 mm glass coverslip for 6-12 hours at 60°C. Arrays are washed for 5 min at 25°C in low stringency wash buffer (1 x SSC/0.2% SDS), then for 10 min at room temperature in high stringency wash buffer (0.1 x SSC/0.2% SDS). Arrays are scanned in 0.1 x SSC using a fluorescence laser scanning device fitted with a custom  
15 filter set. Accurate differential expression measurements are obtained by taking the average of the ratios of two independent hybridizations.

Quantitative analysis of the expression of genes may also be performed with EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments EST-related nucleic acids, or fragments of positional segments of EST-related nucleic acids in complementary DNA arrays as  
20 described by Pietu *et al.* (*Genome Research* 6:492-503, 1996). The EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments EST-related nucleic acids, or fragments of positional segments of EST-related nucleic acids thereof are PCR amplified and spotted on membranes. Then, mRNAs originating from various tissues or cells are labeled with radioactive nucleotides. After hybridization and washing in controlled conditions, the hybridized mRNAs are detected by phospho-  
25 imaging or autoradiography. Duplicate experiments are performed and a quantitative analysis of differentially expressed mRNAs is then performed.

Alternatively, expression analysis of the EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments EST-related nucleic acids, or fragments of positional segments of EST-related nucleic acids can be done through high density nucleotide arrays as described by Lockhart  
30 *et al.* (*Nature Biotechnology* 14: 1675-1680, 1996) and Sosnowsky *et al.* (*Proc. Natl. Acad. Sci.* 94:1119-1123, 1997). Oligonucleotides of 15-50 nucleotides corresponding to sequences of EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments EST-related nucleic acids, or fragments of positional segments of EST-related nucleic acids are synthesized directly on the chip (Lockhart *et al.*, *supra*) or synthesized and then addressed to the chip (Sosnowsky *et al.*, *supra*).  
35 Preferably, the oligonucleotides are about 20 to 25 nucleotides in length.

cDNA probes labeled with an appropriate compound, such as biotin, digoxigenin or fluorescent dye, are synthesized from the appropriate mRNA population and then randomly fragmented to an

average size of 50 to 100 nucleotides. The said probes are then hybridized to the chip. After washing as described in Lockhart *et al, supra* and application of different electric fields (Sonowsky *et al, supra.*), the dyes or labeling compounds are detected and quantified. Duplicate hybridizations are performed. Comparative analysis of the intensity of the signal originating from cDNA probes on the same target  
5 oligonucleotide in different cDNA samples indicates a differential expression of the mRNA corresponding to the 5' EST, consensus contigated 5' EST or extended cDNA from which the oligonucleotide sequence has been designed.

#### IV. Use of 5' ESTs to Clone Extended cDNAs and to Clone the Corresponding Genomic DNAs

10 Once 5' ESTs or consensus contigated 5' ESTs which include the 5' end of the corresponding mRNAs have been selected using the procedures described above, they can be utilized to isolate extended cDNAs which contain sequences adjacent to the 5' ESTs or consensus contigated 5' ESTs. The extended cDNAs may include the entire coding sequence of the protein encoded by the corresponding mRNA, including the authentic translation start site. If the extended cDNA encodes a  
15 secreted protein, it may contain the signal sequence, and the sequence encoding the mature protein remaining after cleavage of the signal peptide.

Extended cDNAs which include the entire coding sequence of the protein encoded by the corresponding mRNA are referred to herein as "full-length cDNAs." Alternatively, the extended cDNAs may not include the entire coding sequence of the protein encoded by the corresponding mRNA,  
20 although they do include sequences adjacent to the 5'ESTs or consensus contigated 5' ESTs. In some embodiments in which the extended cDNAs are derived from an mRNA encoding a secreted protein, the extended cDNAs may include only the sequence encoding the mature protein remaining after cleavage of the signal peptide, or only the sequence encoding the signal peptide.

Examples 19 and 20 below describe a general method for obtaining extended cDNAs using 5'  
25 ESTs or consensus contigated 5' ESTs and nucleic acid homologous thereto. Example 21 below describes the cloning and sequencing of several extended cDNAs, including full-length cDNAs which include the authentic 5' end of the corresponding mRNA for several secreted proteins.

The methods of Examples 19 and 20 can also be used to obtain extended cDNAs which encode less than the entire coding sequence of proteins encoded by the genes corresponding to the 5' ESTs or  
30 consensus contigated 5'ESTs. In some embodiments, the extended cDNAs isolated using these methods encode at least 5, 10, 15, 20, 25, 30, 35, 40, 50, 75, 100, or 150 consecutive amino acids of one of the proteins encoded by the sequences of SEQ ID NOs. 24-811 and 1600-1622. In some embodiments, the extended cDNAs isolated using these methods encode at least 5, 10, 15, 20, 25, 30, 35, 40, 50, 75, 100, or 150 consecutive amino acids of one of the proteins encoded by the sequences of SEQ ID NOs. 24-  
35 811.

#### EXAMPLE 19



General Method for Using 5' ESTs or Consensus Contigated 5'ESTs to Clone and Sequence Extended cDNAs which Include the Entire Coding Region and the Authentic 5'End of the Corresponding mRNA

The following general method may be used to quickly and efficiently isolate extended cDNAs including sequence adjacent to the sequences of the 5' ESTs or Consensus Contigated 5'ESTs used to  
5 obtain them. This method may be applied to obtain extended cDNAs for any 5' EST or consensus contigated 5' EST of the invention, including those 5' ESTs and consensus contigated 5' ESTs encoding secreted proteins. This method is illustrated in Figure 3.

1. Obtaining Extended cDNAs

The method takes advantage of the known 5' sequence of the mRNA. A reverse transcription  
10 reaction is conducted on purified mRNA with a poly dT primer containing a nucleotide sequence at its 5' end allowing the addition of a known sequence at the end of the cDNA which corresponds to the 3' end of the mRNA. Such a primer and a commercially-available reverse transcriptase enzyme are added to a buffered mRNA sample yielding a reverse transcript anchored at the 3' polyA site of the RNAs. Nucleotide monomers are then added to complete the first strand synthesis.

15 After removal of the mRNA hybridized to the first cDNA strand by alkaline hydrolysis, the products of the alkaline hydrolysis and the residual poly dT primer can be eliminated with an exclusion column.

Subsequently, a pair of nested primers on each end is designed based on the known 5' sequence from the 5' EST or consensus contigated 5' EST and the known 3' end added by the poly dT primer  
20 used in the first strand synthesis. Software used to design primers are either based on GC content and melting temperatures of oligonucleotides, such as OSP (Illier and Green, *PCR Meth. Appl.* 1:124-128, 1991), or based on the octamer frequency disparity method (Griffais *et al.*, *Nucleic Acids Res.* 19: 3887-3891, 1991 ) such as PC-Rare ([http:// bioinformatics.weizmann.ac.il/software/PC-Rare/doc/manuel.html](http://bioinformatics.weizmann.ac.il/software/PC-Rare/doc/manuel.html)). Preferably, the nested primers at the 5' end and the nested primers at the 3' end  
25 are separated from one another by four to nine bases. These primer sequences may be selected to have melting temperatures and specificities suitable for use in PCR.

A first PCR run is performed using the outer primer from each of the nested pairs. A second PCR run using the inner primer from each of the nested pairs is then performed on a small sample of the first PCR product. Thereafter, the primers and remaining nucleotide monomers are removed.

30 2. Sequencing Extended cDNAs or Fragments Thereof

Due to the lack of position constraints on the design of 5' nested primers compatible for PCR use using the OSP software, amplicons of two types are obtained. Preferably, the second 5' primer is located upstream of the translation initiation codon thus yielding a nested PCR product containing the entire coding sequence. Such an extended cDNA may be used in a direct cloning procedure as described  
35 in section a below. However, in some cases, the second 5' primer is located downstream of the translation initiation codon, thereby yielding a PCR product containing only part of the ORF. Such incomplete PCR products are submitted to a modified procedure described in section b below.

a) *Nested PCR products containing complete ORFs*

When the resulting nested PCR product contains the complete coding sequence, as predicted from the 5' EST or consensus contigated 5' EST sequence, it is directly cloned in an appropriate vector as described in section 3.

5 b) *Nested PCR products containing incomplete ORFs*

When the amplicon does not contain the complete coding sequence, intermediate steps are necessary to obtain both the complete coding sequence and a PCR product containing the full coding sequence. The complete coding sequence can be assembled from several partial sequences determined directly from different PCR products.

10 Once the full coding sequence has been completely determined, new primers compatible for PCR use are then designed to obtain amplicons containing the whole coding region. However, in such cases, 3' primers compatible for PCR use are located inside the 3' UTR of the corresponding mRNA, thus yielding amplicons which lack part of this region, *i.e.* the polyA tract and sometimes the polyadenylation signal, as illustrated in Figure 3. Such extended cDNAs are then cloned into an  
15 appropriate vector as described in section 3.

c) *Sequencing extended cDNAs*

Sequencing of extended cDNAs can be performed using a Die Terminator approach with the AmpliTaq DNA polymerase FS kit available from Perkin Elmer.

In order to sequence long PCR fragments, primer walking is performed using software such as  
20 OSP to choose primers and automated computer software such as ASMG (Sutton *et al.*, *Genome Science Technol.* 1: 9-19, 1995) to construct contigs of walking sequences including the initial 5' tag. Preferably, primer walking is performed until the sequences of full length cDNAs are obtained.

Completion of the sequencing of a given extended cDNA fragment may be assessed by comparing the sequence length to the size of the corresponding nested PCR product. When Northern  
25 blot data are available, the size of the mRNA detected for a given PCR product may also be used to finally assess that the sequence is complete. Sequences which do not fulfill these criteria are discarded and will undergo a new isolation procedure.

3. Cloning Extended cDNAs

The PCR product containing the full coding sequence is then cloned in an appropriate vector.  
30 For example, the extended cDNAs can be cloned into any expression vector known in the art, such as pED6dpc2 (DiscoverEase, Genetics Institute, Cambridge, MA).

Cloned PCR products are then entirely sequenced in order to obtain at least two sequences per clone. Preferably, the sequences are obtained from both sense and antisense strands according to the aforementioned procedure with the following modifications. First, both 5' and 3' ends of cloned  
35 PCR products are sequenced in order to confirm the identity of the clone. Second, primer walking is performed if the full coding region has not been obtained yet. Contiguation is then performed using primer walking sequences for cloned products as well as walking sequences that have already

contiguated for uncloned PCR products. The sequence is considered complete when the resulting contigs include the whole coding region as well as overlapping sequences with vector DNA on both ends. All the contiguated sequences for each cloned amplicon are then used to obtain a consensus sequence.

#### 5 4. Selection of Cloned Full length Sequences

##### *a) Computer analysis of extended cDNAs*

Following identification of contaminants and masking of repeats, structural features, e.g. polyA tail and polyadenylation signal, of the sequences of extended cDNAs are subsequently determined using methods known to those skilled in the art. For example, algorithm, parameters and  
10 criteria defined in Figure 10 may be used. Briefly, a polyA tail is defined as a homopolymeric stretch of at least 11 A with at most one alternative base within it. The polyA tail search is restricted to the last 20 nucleotides of the sequence and limited to stretches of 11 consecutive A's because sequencing reactions are often not readable after such a polyA stretch. To search for a polyadenylation signal, the polyA tail is clipped from the full-length sequence. The 50 nucleotides preceding the polyA tail  
15 are searched for the canonic polyadenylation AAUAAA signal allowing one mismatch to account for possible sequencing errors as well as known variation in the canonical sequence of the polyadenylation signal.

Functional features, e.g. ORFs and signal sequences, of the sequences of extended cDNAs are subsequently determined as follows. The 3 upper strand frames of extended cDNAs are searched for  
20 ORFs defined as the maximum length fragments beginning with a translation initiation codon and ending with a stop codon. ORFs encoding at least 80 amino acids are preferred. If extended cDNAs encoding secreted proteins are desired, each found ORF is then scanned for the presence of a signal peptide using the matrix method described in Example 13.

Sequences of extended cDNAs are then compared, on a nucleotidic or proteic basis, to public  
25 sequences available at the time of filing.

##### *b) Selection of full-length cDNAs of interest*

A negative selection may then be performed in order to eliminate unwanted cloned sequences resulting from either contaminants or PCR artifacts as follows. Sequences matching contaminant sequences such as vector DNA, tRNA, mtRNA, rRNA sequences are discarded as well as those  
30 encoding ORF sequences exhibiting extensive homology to repeats. Sequences obtained by direct cloning (section 1a) but lacking polyA tail may be discarded. Only ORFs ending either before the polyA tail (section 1a) or before the end of the cloned 3'UTR (section 1b) may be selected. If extended cDNAs encoding secreted proteins are desired, ORFs containing a signal peptide are considered. In addition, ORFs containing unlikely mature proteins such as mature proteins which size is less than 20 amino acids  
35 or less than 25% of the immature protein size may be eliminated.

Then, for each remaining full length cDNA containing several ORFs, a preselection of ORFs may be performed using the following criteria. The longest ORF is preferred. If extended cDNAs

encoding secreted proteins are desired and if the ORF sizes are similar, the chosen ORF is the one which signal peptide has the highest score according to Von Heijne method.

Sequences of full length cDNA clones may then be compared pairwise after masking of the repeat sequences. Full-length cDNA sequences exhibiting extensive homology may be clustered in the same class. Each cluster may then be subjected to a cluster analysis that detects sequences resulting from internal priming or from alternative splicing, identical sequences or sequences with several frameshifts. A selection may be operated between clones belonging to the same class in order to detect clones encoding homologous but distinct ORFs which may be both selected if they both contain sequences of interest.

Selection of full-length cDNA clones encoding sequences of interest may subsequently be performed using the following criteria. Structural parameters (initial tag, polyadenylation site and signal) are first checked. Then, homologies with known nucleic acids and proteins are examined in order to determine whether the clone sequence match a known nucleotide/protein sequence and, in the latter case, its covering rate and the date at which the sequence became public. If there is no extensive match with sequences other than ESTs or genomic DNA, or if the clone sequence brings substantial new information, such as encoding a protein resulting from alternative splicing of an mRNA coding for an already known protein, the sequence is kept. Examples of such cloned full-length cDNAs containing sequences of interest are described in Example 21. Sequences resulting from chimera or double inserts or located on chromosome breaking points as assessed by homology to other sequences may be discarded during this procedure.

Extended cDNAs prepared as described above may be subsequently engineered to obtain nucleic acids which include desired portions of the extended cDNA using conventional techniques such as subcloning, PCR, or *in vitro* oligonucleotide synthesis. For example, nucleic acids which include only the full coding sequences may be obtained using techniques known to those skilled in the art.

Alternatively, conventional techniques may be applied to obtain nucleic acids which contain only part of the coding sequences. In the case of nucleic acids encoding secreted proteins, nucleic acids containing only the coding sequence for the mature protein remaining after the signal peptide is cleaved off or nucleic acids which contain only the coding sequences for the signal peptides may be obtained.

Similarly, nucleic acids containing any other desired portion of the coding sequences for the encoded protein may be obtained. For example, the nucleic acid may contain at least 10, 15, 18, 20, 25, 28, 30, 35, 40, 50, 75, 100, 150, 200, 300, 400 or 500 consecutive bases of an extended cDNA.

Once an extended cDNA has been obtained, it can be sequenced to determine the amino acid sequence it encodes. Once the encoded amino acid sequence has been determined, one can create and identify any of the many conceivable cDNAs that will encode that protein by simply using the degeneracy of the genetic code. For example, allelic variants or other homologous nucleic acids can be identified as described below. Alternatively, nucleic acids encoding the desired amino acid sequence can be synthesized *in vitro*.

In a preferred embodiment, the coding sequence may be selected using the known codon or codon pair preferences for the host organism in which the cDNA is to be expressed.

In addition to PCR based methods for obtaining cDNAs which include the authentic 5' end of the corresponding mRNA as well as the complete protein coding sequence of the corresponding mRNA, traditional hybridization based methods may also be employed. These methods may also be used to obtain the genomic DNAs which encode the mRNAs from which the 5' ESTs or consensus contigated 5' ESTS were derived, mRNAs corresponding to the extended cDNAs, or nucleic acids which are homologous to extended cDNAs, 5' ESTs, or consensus contigated 5' ESTs. Example 19 below provides examples of such methods.

#### EXAMPLE 20

##### Methods for Obtaining Extended cDNAs which Include the Entire Coding Region and the Authentic 5' End of the Corresponding mRNA or Nucleic Acids Homologous to Extended cDNAs, 5' ESTs or Consensus Contigated 5' ESTs

A full-length cDNA library can be made using the strategies described in Example 7. Alternatively, a cDNA library or genomic DNA library may be obtained from a commercial source or made using techniques familiar to those skilled in the art.

Such cDNA or genomic DNA libraries may be used to isolate extended cDNAs obtained from 5' ESTs or consensus contigated 5' ESTs or nucleic acids homologous to extended cDNAs, 5' ESTs, or consensus contigated 5' ESTs as follows. The cDNA library or genomic DNA library is hybridized to a detectable probe. The detectable probe may comprise at least 10, 15, 18, 20, 25, 28, 30, 35, 40, 50, 75, 100, 150, 200, 300, 400 or 500 consecutive nucleotides of the 5' EST, consensus contigated 5' EST, or extended cDNA.

Techniques for identifying cDNA clones in a cDNA library which hybridize to a given probe sequence are disclosed in Sambrook *et al.*, *Molecular Cloning: A Laboratory Manual 2d Ed.*, Cold Spring Harbor Laboratory Press, 1989. The same techniques may be used to isolate genomic DNAs. Briefly, cDNA or genomic DNA clones which hybridize to the detectable probe are identified and isolated for further manipulation as follows. The detectable probe described in the preceding paragraph is labeled with a detectable label such as a radioisotope or a fluorescent molecule. Techniques for labeling the probe are well known and include phosphorylation with polynucleotide kinase, nick translation, *in vitro* transcription, and non radioactive techniques. The cDNAs or genomic DNAs in the library are transferred to a nitrocellulose or nylon filter and denatured. After blocking of non specific sites, the filter is incubated with the labeled probe for an amount of time sufficient to allow binding of the probe to cDNAs or genomic DNAs containing a sequence capable of hybridizing thereto.

By varying the stringency of the hybridization conditions used to identify cDNAs or genomic DNAs which hybridize to the detectable probe, cDNAs or genomic DNAs having different levels of homology to the probe can be identified and isolated as described below.

### 1. Identification of cDNA or Genomic DNA Sequences Having a High Degree of Homology to the Labeled Probe

To identify cDNAs or genomic DNAs having a high degree of homology to the probe sequence, the melting temperature of the probe may be calculated using the following formulas:

- 5 For probes between 14 and 70 nucleotides in length the melting temperature ( $T_m$ ) is calculated using the formula:  $T_m = 81.5 + 16.6(\log(\text{Na}^+)) + 0.41(\text{fraction G+C}) - (600/N)$  where N is the length of the probe.

If the hybridization is carried out in a solution containing formamide, the melting temperature may be calculated using the equation  $T_m = 81.5 + 16.6(\log(\text{Na}^+)) + 0.41(\text{fraction G+C}) - (0.63\%$

- 10 formamide)  $- (600/N)$  where N is the length of the probe.

Prehybridization may be carried out in 6X SSC, 5X Denhardt's reagent, 0.5% SDS, 100  $\mu\text{g}$  denatured fragmented salmon sperm DNA or 6X SSC, 5X Denhardt's reagent, 0.5% SDS, 100  $\mu\text{g}$  denatured fragmented salmon sperm DNA, 50% formamide. The formulas for SSC and Denhardt's solutions are listed in Sambrook *et al.*, *supra*.

- 15 Hybridization is conducted by adding the detectable probe to the prehybridization solutions listed above. Where the probe comprises double stranded DNA, it is denatured before addition to the hybridization solution. The filter is contacted with the hybridization solution for a sufficient period of time to allow the probe to hybridize to extended cDNAs or genomic DNAs containing sequences complementary thereto or homologous thereto. For probes over 200 nucleotides in length, the
- 20 hybridization may be carried out at 15-25°C below the  $T_m$ . For shorter probes, such as oligonucleotide probes, the hybridization may be conducted at 15-25°C below the  $T_m$ . Preferably, for hybridizations in 6X SSC, the hybridization is conducted at approximately 68°C. Preferably, for hybridizations in 50% formamide containing solutions, the hybridization is conducted at approximately 42°C.

All of the foregoing hybridizations would be considered to be under "stringent" conditions.

- 25 Following hybridization, the filter is washed in 2X SSC, 0.1% SDS at room temperature for 15 minutes. The filter is then washed with 0.1X SSC, 0.5% SDS at room temperature for 30 minutes to 1 hour. Thereafter, the solution is washed at the hybridization temperature in 0.1X SSC, 0.5% SDS. A final wash is conducted in 0.1X SSC at room temperature.

- cDNAs or genomic DNAs which have hybridized to the probe are identified by autoradiography
- 30 or other conventional techniques.

### 2. Obtaining cDNA or Genomic DNA Sequences Having Lower Degrees of Homology to the Labeled Probe

The above procedure may be modified to identify cDNAs or genomic DNAs having decreasing levels of homology to the probe sequence. For example, to obtain cDNAs or genomic DNAs of

- 35 decreasing homology to the detectable probe, less stringent conditions may be used. For example, the hybridization temperature may be decreased in increments of 5°C from 68°C to 42°C in a hybridization buffer having a sodium concentration of approximately 1M. Following hybridization, the filter may be

washed with 2X SSC, 0.5% SDS at the temperature of hybridization. These conditions are considered to be "moderate" conditions above 50°C and "low" conditions below 50°C.

Alternatively, the hybridization may be carried out in buffers, such as 6X SSC, containing formamide at a temperature of 42°C. In this case, the concentration of formamide in the hybridization buffer may be reduced in 5% increments from 50% to 0% to identify clones having decreasing levels of homology to the probe. Following hybridization, the filter may be washed with 6X SSC, 0.5% SDS at 50°C. These conditions are considered to be "moderate" conditions above 25% formamide and "low" conditions below 25% formamide. cDNAs or genomic DNAs which have hybridized to the probe are identified by autoradiography.

10 3. Determination of the Degree of Homology between the Obtained cDNAs or Genomic DNAs and 5'ESTs, Consensus Contigated 5'ESTs, or Extended cDNAs or Between the Polypeptides Encoded by the Obtained cDNAs or Genomic DNAs and the Polypeptides Encoded by the 5'ESTs, Consensus Contigated 5'ESTs, or Extended cDNAs

To determine the level of homology between the hybridized cDNA or genomic DNA and the 5'EST, consensus contigated 5'EST or extended cDNA from which the probe was derived, the nucleotide sequences of the hybridized nucleic acid and the 5'EST, consensus contigated 5'EST or extended cDNA from which the probe was derived are compared. The sequences of the 5'EST, consensus contigated 5'EST or extended cDNA from which the probe was derived and the sequences of the cDNA or genomic DNA which hybridized to the detectable probe may be stored on a computer readable medium as described below and compared to one another using any of a variety of algorithms familiar to those skilled in the art, those described below.

To determine the level of homology between the polypeptide encoded by the hybridizing cDNA or genomic DNA and the polypeptide encoded by the 5'EST, consensus contigated 5'EST or extended cDNA from which the probe was derived, the polypeptide sequence encoded by the hybridized nucleic acid and the polypeptide sequence encoded by the 5'EST, consensus contigated 5'EST or extended cDNA from which the probe was derived are compared. The sequences of the polypeptide encoded by the 5'EST, consensus contigated 5'EST or extended cDNA from which the probe was derived and the polypeptide sequence encoded by the cDNA or genomic DNA which hybridized to the detectable probe may be stored on a computer readable medium as described below and compared to one another using any of a variety of algorithms familiar to those skilled in the art, those described below.

Protein and/or nucleic acid sequence homologies may be evaluated using any of the variety of sequence comparison algorithms and programs known in the art. Such algorithms and programs include, but are by no means limited to, TBLASTN, BLASTP, FASTA, TFASTA, and CLUSTALW (Pearson and Lipman, 1988, *Proc. Natl. Acad. Sci. USA* 85(8):2444-2448; Altschul *et al.*, 1990, *J. Mol. Biol.* 215(3):403-410; Thompson *et al.*, 1994, *Nucleic Acids Res.* 22(2):4673-4680; Higgins *et al.*, 1996, *Methods Enzymol.* 266:383-402; Altschul *et al.*, 1990, *J. Mol. Biol.* 215(3):403-410; Altschul *et al.*, 1993, *Nature Genetics* 3:266-272).

In a particularly preferred embodiment, protein and nucleic acid sequence homologies are evaluated using the Basic Local Alignment Search Tool ("BLAST") which is well known in the art (see, e.g., Karlin and Altschul, 1990, *Proc. Natl. Acad. Sci. USA* 87:2267-2268; Altschul *et al.*, 1990, *J. Mol. Biol.* 215:403-410; Altschul *et al.*, 1993, *Nature Genetics* 3:266-272; Altschul *et al.*, 1997, *Nuc. Acids Res.* 25:3389-3402). In particular, five specific BLAST programs are used to perform the following task:

- (1) BLASTP and BLAST3 compare an amino acid query sequence against a protein sequence database;
- (2) BLASTN compares a nucleotide query sequence against a nucleotide sequence database;
- (3) BLASTX compares the six-frame conceptual translation products of a query nucleotide sequence (both strands) against a protein sequence database;
- (4) TBLASTN compares a query protein sequence against a nucleotide sequence database translated in all six reading frames (both strands); and
- 15 (5) TBLASTX compares the six-frame translations of a nucleotide query sequence against the six-frame translations of a nucleotide sequence database.

The BLAST programs identify homologous sequences by identifying similar segments, which are referred to herein as "high-scoring segment pairs," between a query amino or nucleic acid sequence and a test sequence which is preferably obtained from a protein or nucleic acid sequence database. High-scoring segment pairs are preferably identified (*i.e.*, aligned) by means of a scoring matrix, many of which are known in the art. Preferably, the scoring matrix used is the BLOSUM62 matrix (Gonnet *et al.*, 1992, *Science* 256:1443-1445; Henikoff and Henikoff, 1993, *Proteins* 17:49-61). Less preferably, the PAM or PAM250 matrices may also be used (see, e.g., Schwartz and Dayhoff, eds., 1978, *Matrices for Detecting Distance Relationships: Atlas of Protein Sequence and*  
25 *Structure*, Washington: National Biomedical Research Foundation)

The BLAST programs evaluate the statistical significance of all high-scoring segment pairs identified, and preferably selects those segments which satisfy a user-specified threshold of significance, such as a user-specified percent homology. Preferably, the statistical significance of a high-scoring segment pair is evaluated using the statistical significance formula of Karlin (see, e.g.,  
30 Karlin and Altschul, 1990, *Proc. Natl. Acad. Sci. USA* 87:2267-2268).

The parameters used with the above algorithms may be adapted depending on the sequence length and degree of homology studied. In some embodiments, the parameters may be the default parameters used by the algorithms in the absence of instructions from the user.

In some embodiments, the level of homology between the hybridized nucleic acid and the  
35 extended cDNA, 5'EST, or 5' consensus contigated 5'EST from which the probe was derived may be determined using the FASTDB algorithm described in Brutlag *et al.* *Comp. App. Biosci.* 6:237-245, 1990. In such analyses the parameters may be selected as follows: Matrix=Unitary, k-tuple=4,



Mismatch Penalty=1, Joining Penalty=30, Randomization Group Length=0, Cutoff Score=1, Gap Penalty=5, Gap Size Penalty=0.05, Window Size=500 or the length of the sequence which hybridizes to the probe, whichever is shorter. Because the FASTDB program does not consider 5' or 3' truncations when calculating homology levels, if the sequence which hybridizes to the probe is truncated relative to the sequence of the extended cDNA, 5'EST, or consensus contigated 5'EST from which the probe was derived the homology level is manually adjusted by calculating the number of nucleotides of the extended cDNA, 5'EST, or consensus contigated 5' EST which are not matched or aligned with the hybridizing sequence, determining the percentage of total nucleotides of the hybridizing sequence which the non-matched or non-aligned nucleotides represent, and subtracting this percentage from the homology level. For example, if the hybridizing sequence is 700 nucleotides in length and the extended cDNA, 5'EST, or consensus contigated 5' EST sequence is 1000 nucleotides in length wherein the first 300 bases at the 5' end of the extended cDNA, 5'EST, or consensus contigated 5' EST are absent from the hybridizing sequence, and wherein the overlapping 700 nucleotides are identical, the homology level would be adjusted as follows. The non-matched, non-aligned 300 bases represent 30% of the length of the extended cDNA, 5'EST, or consensus contigated 5' EST. If the overlapping 700 nucleotides are 100% identical, the adjusted homology level would be  $100-30=70\%$  homology. It should be noted that the preceding adjustments are only made when the non-matched or non-aligned nucleotides are at the 5' or 3' ends. No adjustments are made if the non-matched or non-aligned sequences are internal or under any other conditions.

For example, using the above methods, nucleic acids having at least 95% nucleic acid homology, at least 96% nucleic acid homology, at least 97% nucleic acid homology, at least 98% nucleic acid homology, at least 99% nucleic acid homology, or more than 99% nucleic acid homology to the extended cDNA, 5'EST, or consensus contigated 5' EST from which the probe was derived may be obtained and identified. Such nucleic acids may be allelic variants or related nucleic acids from other species. Similarly, by using progressively less stringent hybridization conditions one can obtain and identify nucleic acids having at least 90%, at least 85%, at least 80% or at least 75% homology to the extended cDNA, 5'EST, or consensus contigated 5' EST from which the probe was derived.

Using the above methods and algorithms such as FASTA with parameters depending on the sequence length and degree of homology studied, for example the default parameters used by the algorithms in the absence of instructions from the user, one can obtain nucleic acids encoding proteins having at least 99%, at least 98%, at least 97%, at least 96%, at least 95%, at least 90%, at least 85%, at least 80% or at least 75% homology to the protein encoded by the extended cDNA, 5'EST, or consensus contigated 5' EST from which the probe was derived. In some embodiments, the homology levels can be determined using the "default" opening penalty and the "default" gap penalty, and a scoring matrix such as PAM 250 (a standard scoring matrix; see Dayhoff *et al.*, in: Atlas of Protein Sequence and Structure, Vol. 5, Supp. 3 (1978)).

Alternatively, the level of polypeptide homology may be determined using the FASTDB algorithm described by Brutlag *et al.* Comp. App. Biosci. 6:237-245, 1990. In such analyses the parameters may be selected as follows: Matrix=PAM 0, k-tuple=2, Mismatch Penalty=1, Joining Penalty=20, Randomization Group Length=0, Cutoff Score=1, Window Size=Sequence Length, Gap Penalty=5, Gap Size Penalty=0.05, Window Size=500 or the length of the homologous sequence, whichever is shorter. If the homologous amino acid sequence is shorter than the amino acid sequence encoded by the extended cDNA, 5'EST, or consensus contigated 5' EST as a result of an N terminal and/or C terminal deletion the results may be manually corrected as follows. First, the number of amino acid residues of the amino acid sequence encoded by the extended cDNA, 5'EST, or consensus contigated 5' EST which are not matched or aligned with the homologous sequence is determined. Then, the percentage of the length of the sequence encoded by the extended cDNA, 5'EST, or consensus contigated 5' EST which the non-matched or non-aligned amino acids represent is calculated. This percentage is subtracted from the homology level. For example wherein the amino acid sequence encoded by the extended cDNA, 5'EST, or consensus contigated 5' EST is 100 amino acids in length and the length of the homologous sequence is 80 amino acids and wherein the amino acid sequence encoded by the extended cDNA or 5'EST is truncated at the N terminal end with respect to the homologous sequence, the homology level is calculated as follows. In the preceding scenario there are 20 non-matched, non-aligned amino acids in the sequence encoded by the extended cDNA, 5'EST, or consensus contigated 5' EST. This represents 20% of the length of the amino acid sequence encoded by the extended cDNA, 5'EST, or consensus contigated 5' EST. If the remaining amino acids are 100% identical between the two sequences, the homology level would be  $100\% - 20\% = 80\%$  homology. No adjustments are made if the non-matched or non-aligned sequences are internal or under any other conditions.

In addition to the above described methods, other protocols are available to obtain extended cDNAs using 5' ESTs or consensus contigated 5'ESTs as outlined in the following paragraphs.

Extended cDNAs may be prepared by obtaining mRNA from the tissue, cell, or organism of interest using mRNA preparation procedures utilizing polyA selection procedures or other techniques known to those skilled in the art. A first primer capable of hybridizing to the polyA tail of the mRNA is hybridized to the mRNA and a reverse transcription reaction is performed to generate a first cDNA strand.

The first cDNA strand is hybridized to a second primer containing at least 10 consecutive nucleotides of the sequences of SEQ ID NOs 24-811 and 1600-1622. Preferably, the primer comprises at least 10, 12, 15, 17, 18, 20, 23, 25, or 28 consecutive nucleotides from the sequences of SEQ ID NOs 24-811 and 1600-1622. In some embodiments, the primer comprises more than 30 nucleotides from the sequences of SEQ ID NOs 24-811 and 1600-1622. If it is desired to obtain extended cDNAs containing the full protein coding sequence, including the authentic translation initiation site, the second primer used contains sequences located upstream of the translation initiation site. The second primer is

extended to generate a second cDNA strand complementary to the first cDNA strand. Alternatively, RT-PCR may be performed as described above using primers from both ends of the cDNA to be obtained.

Extended cDNAs containing 5' fragments of the mRNA may be prepared by hybridizing an mRNA comprising the sequences of SEQ ID NOs. 24-811 and 1600-1622 with a primer comprising a  
5 complementary to a fragment of an EST-related nucleic acid hybridizing the primer to the mRNAs, and reverse transcribing the hybridized primer to make a first cDNA strand from the mRNAs. Preferably, the primer comprises at least 10, 12, 15, 17, 18, 20, 23, 25, or 28 consecutive nucleotides of the sequences complementary to SEQ ID NOs. 24-811 and 1600-1622.

Thereafter, a second cDNA strand complementary to the first cDNA strand is synthesized. The  
10 second cDNA strand may be made by hybridizing a primer complementary to sequences in the first cDNA strand to the first cDNA strand and extending the primer to generate the second cDNA strand.

The double stranded extended cDNAs made using the methods described above are isolated and cloned. The extended cDNAs may be cloned into vectors such as plasmids or viral vectors capable of replicating in an appropriate host cell. For example, the host cell may be a bacterial, mammalian,  
15 avian, or insect cell.

Techniques for isolating mRNA, reverse transcribing a primer hybridized to mRNA to generate a first cDNA strand, extending a primer to make a second cDNA strand complementary to the first cDNA strand, isolating the double stranded cDNA and cloning the double stranded cDNA are well known to those skilled in the art and are described in *Current Protocols in Molecular Biology*, John  
20 Wiley & Sons, Inc. 1997 and Sambrook *et al.*, *Molecular Cloning: A Laboratory Manual*, Second Edition, Cold Spring Harbor Laboratory Press, 1989.

Alternatively, other procedures may be used for obtaining full-length cDNAs or extended cDNAs. In one approach, full-length or extended cDNAs are prepared from mRNA and cloned into double stranded phagemids as follows. The cDNA library in the double stranded phagemids is then  
25 rendered single stranded by treatment with an endonuclease, such as the Gene II product of the phage F1 and an exonuclease (Chang *et al.*, *Gene* 127:95-8, 1993). A biotinylated oligonucleotide comprising the sequence of a fragment of an EST-related nucleic acid is hybridized to the single stranded phagemids. Preferably, the fragment comprises at least 10, 12, 15, 17, 18, 20, 23, 25, or 28 consecutive nucleotides of the sequences of SEQ ID NOs. 24-811 and 1600-1622.

30 Hybrids between the biotinylated oligonucleotide and phagemids are isolated by incubating the hybrids with streptavidin coated paramagnetic beads and retrieving the beads with a magnet (Fry *et al.*, *Biotechniques*, 13: 124-131, 1992). Thereafter, the resulting phagemids are released from the beads and converted into double stranded DNA using a primer specific for the 5' EST or consensus contigated 5'EST sequence used to design the biotinylated oligonucleotide. Alternatively, protocols such as the  
35 Gene Trapper kit (Gibco BRL) may be used. The resulting double stranded DNA is transformed into bacteria. Extended cDNAs or full length cDNAs containing the 5' EST or consensus contigated 5'EST sequence are identified by colony PCR or colony hybridization.

Using any of the above described methods in section III, a plurality of extended cDNAs containing full-length protein coding sequences or portions of the protein coding sequences may be provided as cDNA libraries for subsequent evaluation of the encoded proteins or use in diagnostic assays as described below.

5

### EXAMPLE 21

#### Full Length cDNAs

The procedures described in Example 19 and 20 were used to obtain extended cDNAs or full length cDNAs derived from 5' ESTs in a variety of tissues. The following list provides a few examples  
10 of cDNAs obtained by these means.

Using this procedure, the full length cDNA of SEQ ID NO:1 (internal identification number 58-34-2-E7-FL2) was obtained. This cDNA encodes the signal peptide MWWFQQGLSFLPSALVIWTSA (SEQ ID NO:2) having a von Heijne score of 5.5.

Using this approach, the full length cDNA of SEQ ID NO:3 (internal identification number 48-  
15 19-3-G1-FL1) was obtained. This cDNA encodes the signal peptide MKKVLLITAILAVAVG (SEQ ID NO: 4) having a von Heijne score of 8.2.

The full length cDNA of SEQ ID NO:5 (internal identification number 58-35-2-F10-FL2) was also obtained using this procedure. This cDNA encodes a signal peptide LWLLFFLVTAIHA (SEQ ID NO:6) having a von Heijne score of 10.7.

Furthermore, the polypeptides encoded by the extended or full-length cDNAs may be screened  
20 for the presence of known structural or functional motifs or for the presence of signatures, small amino acid sequences which are well conserved amongst the members of a protein family. The results obtained for the polypeptides encoded by a few full-length cDNAs derived from 5'ESTs that were screened for the presence of known protein signatures and motifs using the Proscan software from the GCG  
25 package and the Prosite 15.0 database are provided below.

The protein of SEQ ID NO: 8 encoded by the full-length cDNA SEQ ID NO: 7 (internal designation 78-8-3-E6-CL0\_1C) and expressed in adult prostate belong to the phosphatidylethanolamine-binding protein from which it exhibits the characteristic PROSITE signature from positions 90 to 112. Proteins from this widespread family, from nematodes to fly,  
30 yeast, rodent and primate species, bind hydrophobic ligands such as phospholipids and nucleotides. They are mostly expressed in brain and in testis and are thought to play a role in cell growth and/or maturation, in regulation of the sperm maturation, motility and in membrane remodeling. They may act either through signal transduction or through oxidoreduction reactions (for a review see Schoentgen and Jollès, *FEBS Letters*, 369 :22-26 (1995)). Taken together, these data suggest that the  
35 protein of SEQ ID NO: 8 may play a role in cell growth, maturation and in membrane remodeling and/or may be related to male fertility. Thus, these protein may be useful in diagnosing and/or treating cancer, neurodegenerative diseases, and/or disorders related to male fertility and sterility.

The protein of SEQ ID No. 10 encoded by the full-length cDNA SEQ ID NO. 9 (internal designation 108-013-5-O-H9-FLC) shows homologies with a family of lysophospholipases conserved among eukaryotes (yeast, rabbit, rodents and human). In addition, some members of this family exhibit a calcium-independent phospholipase A2 activity (Portilla *et al*, *J. Am. Soc. Nephro.*, 9 :1178-1186 (1998)). All members of this family exhibit the active site consensus GX SXG motif of carboxylesterases that is also found in the protein of SEQ ID NO. 10 (position 54 to 58). In addition, this protein may be a membrane protein with one transmembrane domain as predicted by the software TopPred II (Claros and von Heijne, *CABIOS applic. Notes*, 10 :685-686 (1994)). Taken together, these data suggest that the protein of SEQ ID NO:10 may play a role in fatty acid metabolism, probably as a phospholipase. Thus, this protein or part therein, may be useful in diagnosing and/or treating several disorders including, but not limited to, cancer, diabetes, and neurodegenerative disorders such as Parkinson's and Alzheimer's diseases. It may also be useful in modulating inflammatory responses to infectious agents and/or to suppress graft rejection.

The protein of SEQ ID NO: 12 encoded by the full-length cDNA SEQ ID NO: 11 (internal designation 108-004-5-0-D10-FLC) shows remote homology to a subfamily of beta4-galactosyltransferases widely conserved in animals (human, rodents, cow and chicken). Such enzymes, usually type II membrane proteins located in the endoplasmic reticulum or in the Golgi apparatus, catalyzes the biosynthesis of glycoproteins, glycolipid glycans and lactose. Their characteristic features defined as those of subfamily A in Breton *et al*, *J. Biochem.*, 123:1000-1009 (1998) are pretty well conserved in the protein of SEQ ID NO: 12, especially the region I containing the DVD motif (positions 163-165) thought to be involved either in UDP binding or in the catalytic process itself. In addition, the protein of SEQ ID NO: 12 has the typical structure of a type II protein. Indeed, it contains a short 28-amino-acid-long N-terminal tail, a transmembrane segment from positions 29 to 49 and a large 278-amino-acid-long C-terminal tail as predicted by the software TopPred II (Claros and von Heijne, *CABIOS applic. Notes*, 10 :685-686 (1994)). Taken together, these data suggest that the protein of SEQ ID NO: 12 may play a role in the biosynthesis of polysaccharides, and of the carbohydrate moieties of glycoproteins and glycolipids and/or in cell-cell recognition. Thus, this protein may be useful in diagnosing and/or treating several types of disorders including, but not limited to, cancer, atherosclerosis, cardiovascular disorders, autoimmune disorders and rheumatic diseases including rheumatoid arthritis.

The protein of SEQ ID NO: 14 encoded by the full-length cDNA SEQ ID NO: 13 (internal designation 108-009-5-0-A2-FLC) shows extensive homology to the bZIP family of transcription factors, and especially to the human protein (Lu *et al.*, *Mol. Cell. Biol.*, 17 :5117-5126 (1997)). The match include the whole bZIP domain composed of a basic DNA-binding domain and of a leucine zipper allowing protein dimerization. The basic domain is conserved in the protein of SEQ ID NO: 14 as shown by the characteristic PROSITE signature (positions 224-237) except for a conservative substitution of a glutamic acid with an aspartic acid in position 233. The typical

PROSITE signature for leucine zipper is also present (positions 259 to 280). Taken together, these data suggest that the protein of SEQ ID NO: 14 may bind to DNA, hence regulating gene expression as a transcription factor. Thus, this protein may be useful in diagnosing and/or treating several types of disorders including, but not limited to, cancer.

5 Bacterial clones containing plasmids containing the full length cDNAs described above are presently stored in the inventor's laboratories under the internal identification numbers provided above. The inserts may be recovered from the deposited materials by growing an aliquot of the appropriate bacterial clone in the appropriate medium. The plasmid DNA can then be isolated using plasmid isolation procedures familiar to those skilled in the art such as alkaline lysis minipreps or large scale  
10 alkaline lysis plasmid isolation procedures. If desired the plasmid DNA may be further enriched by centrifugation on a cesium chloride gradient, size exclusion chromatography, or anion exchange chromatography. The plasmid DNA obtained using these procedures may then be manipulated using standard cloning techniques familiar to those skilled in the art. Alternatively, a PCR can be done with primers designed at both ends of the insertion. The PCR product which corresponds to the cDNA insert  
15 can then be manipulated using standard cloning techniques familiar to those skilled in the art.

#### **V. Expression of Proteins or Polypeptides Encoded by EST-related nucleic acids or Fragments thereof**

EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of EST-  
20 related nucleic acids, and fragments of positional segments of EST-related nucleic acids may be used to express the polypeptides which they encode. In particular, they may be used to express EST-related polypeptides, fragments of EST-related polypeptides, positional segments of EST-related polypeptides, or fragments of positional segments of EST-related polypeptides. In some embodiments, the EST-related nucleic acids, positional segments of EST-related nucleic acids, and fragments of positional  
25 segments of EST-related nucleic acids may be used to express the full polypeptide (*i.e.* the signal peptide and the mature polypeptide) of a secreted protein, the mature protein (*i.e.* the polypeptide generated after cleavage of the signal peptide), or the signal peptide of a secreted protein. If desired, nucleic acids encoding the signal peptide may be used to facilitate secretion of the expressed protein. It will be appreciated that a plurality of EST-related nucleic acids, fragments of EST-related nucleic acids,  
30 positional segments of EST-related nucleic acids, or fragments of positional segments of EST-related nucleic acids may be simultaneously cloned into expression vectors to create an expression library for analysis of the encoded proteins as described below.

#### **EXAMPLE 22**

35 Expression of the Proteins Encoded by the Genes Corresponding to the  
5'ESTs or Consensus Contigated 5' ESTs

To express their encoded proteins, the EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of EST-related nucleic acids, or fragments of positional segments of EST-related nucleic acids are cloned into a suitable expression vector. In some instances, nucleic acids encoding EST-related polypeptides, fragments of EST-related polypeptides, positional segments of EST-related polypeptides or fragments of positional segments of EST-related polypeptides may be cloned into a suitable expression vector.

In some embodiments, the nucleic acids inserted into the expression vector may comprise the coding sequence of a sequence selected from the group consisting of SEQ ID NOs. 24-811. In other embodiments, the nucleic acids inserted into the expression vector may comprise the full coding sequence (*i.e.* the nucleotides encoding the signal peptide and the mature polypeptide) of one of SEQ ID Nos. 766-792. In some embodiments, the nucleic acid inserted into the expression vector may comprise the nucleotides of one of the sequences of SEQ ID Nos. 766-792 which encode the mature polypeptide (*i.e.* the nucleotides encoding the polypeptide generated after cleavage of the signal peptide). In further embodiments, the nucleic acids inserted into the expression vector may comprise the nucleotides of 24-728 and 766-792 which encode the signal peptide to facilitate secretion of the expressed protein. The nucleic acids inserted into the expression vectors may also contain sequences upstream of the sequences encoding the signal peptide, such as sequences which regulate expression levels or sequences which confer tissue specific expression.

The nucleic acid inserted into the expression vector may encode a polypeptide comprising the one of the sequences of SEQ ID Nos. 812-1599. In some embodiments, the nucleic acid inserted into the expression vector may encode the full polypeptide sequence (*i.e.* the signal peptide and the mature polypeptide) included in one of SEQ ID Nos. 1554-1580. In other embodiments, the nucleic acid inserted into the expression vector may encode the mature polypeptide (*i.e.* the polypeptide generated after cleavage of the signal peptide) included in one of the sequences of SEQ ID Nos. 1554-1580. In further embodiments, the nucleic acids inserted into the expression vector may encode the signal peptide included in one of the sequences of 812-1516 and 1554-1580.

The nucleic acid encoding the protein or polypeptide to be expressed is operably linked to a promoter in an expression vector using conventional cloning technology. The expression vector may be any of the mammalian, yeast, insect or bacterial expression systems known in the art. Commercially available vectors and expression systems are available from a variety of suppliers including Genetics Institute (Cambridge, MA), Stratagene (La Jolla, California), Promega (Madison, Wisconsin), and Invitrogen (San Diego, California). If desired, to enhance expression and facilitate proper protein folding, the codon context and codon pairing of the sequence may be optimized for the particular expression organism in which the expression vector is introduced, as explained by Hatfield, *et al.*, U.S. Patent No. 5,082,767.

The following is provided as one exemplary method to express the proteins encoded by the nucleic acids described above. In some instances the nucleic acid encoding the protein or polypeptide to

be expressed includes a methionine initiation codon and a polyA signal. If the nucleic acid encoding the polypeptide to be expressed lacks a methionine to serve as the initiation site, an initiating methionine can be introduced next to the first codon of the nucleic acid using conventional techniques. Similarly, if the nucleic acid encoding the protein or polypeptide to be expressed lacks a polyA signal, this sequence can  
5 be added to the construct by, for example, splicing out the polyA signal from pSG5 (Stratagene) using BglII and SalI restriction endonuclease enzymes and incorporating it into the mammalian expression vector pXT1 (Stratagene). pXT1 contains the LTRs and a portion of the *gag* gene from Moloney Murine Leukemia Virus. The position of the LTRs in the construct allow efficient stable transfection. The vector includes the Herpes Simplex thymidine kinase promoter and the selectable neomycin gene.  
10 The nucleic acid encoding the polypeptide to be expressed is obtained by PCR from the bacterial vector using oligonucleotide primers complementary to the nucleic acid encoding the protein or polypeptide to be expressed and containing restriction endonuclease sequences for Pst I incorporated into the 5' primer and BglII at the 5' end of 3' primer, taking care to ensure that the nucleic acid encoding the protein or polypeptide to be expressed is correctly positioned with respect to the poly A signal. The purified  
15 fragment obtained from the resulting PCR reaction is digested with PstI, blunt ended with an exonuclease, digested with Bgl II, purified and ligated to pXT1, now containing a poly A signal and digested with BglII.

The ligated product is transfected into mouse NIH 3T3 cells using Lipofectin (Life Technologies, Inc., Grand Island, New York) under conditions outlined in the product specification.  
20 Positive transfectants are selected after growing the transfected cells in 600 µg/ml G418 (Sigma, St. Louis, Missouri).

Alternatively, the nucleic acid encoding the protein or polypeptide to be expressed may be cloned into pED6dpc2. The resulting pED6dpc2 constructs may be transfected into a suitable host cell, such as COS 1 cells. Methotrexate resistant cells are selected and expanded. The expressed protein or  
25 polypeptide may be isolated, purified, or enriched as described above.

To confirm expression of the desired protein or polypeptide, the proteins or polypeptides produced by cells containing a vector with a nucleic acid insert encoding the protein or polypeptide are compared to those lacking such an insert. The expressed proteins are detected using techniques familiar to those skilled in the art such as Coomassie blue or silver staining or using antibodies against the protein  
30 or polypeptide encoded by the nucleic acid insert. Antibodies capable of specifically recognizing the protein of interest may be generated using synthetic 15-mer peptides having a sequence encoded by the appropriate nucleic acid. The synthetic peptides are injected into mice to generate antibody to the polypeptide encoded by the nucleic acid.

If the proteins or polypeptides encoded by the nucleic acid inserts are secreted, medium  
35 prepared from the host cells or organisms containing an expression vector which contains a nucleic acid insert encoding the desired protein or polypeptide is compared to medium prepared from the control cells or organism. The presence of a band in medium from the cells containing the nucleic acid insert which



is absent from preparations from the control cells indicates that the protein or polypeptide encoded by the nucleic acid insert is being expressed and secreted. Generally, the band corresponding to the protein encoded by the nucleic acid insert will have a mobility near that expected based on the number of amino acids in the open reading frame of the nucleic acid insert. However, the band may have a mobility  
5 different than that expected as a result of modifications such as glycosylation, ubiquitination, or enzymatic cleavage.

Alternatively, if the protein expressed from the above expression vectors does not contain sequences directing its secretion, the proteins expressed from host cells containing an expression vector with an insert encoding a secreted protein or portion thereof can be compared to the proteins expressed  
10 in control host cells containing the expression vector without an insert. The presence of a band in samples from cells containing the expression vector with an insert which is absent in samples from cells containing the expression vector without an insert indicates that the desired protein or portion thereof is being expressed. Generally, the band will have the mobility expected for the secreted protein or portion thereof. However, the band may have a mobility different than that expected as a result of modifications  
15 such as glycosylation, ubiquitination, or enzymatic cleavage.

The expressed protein or polypeptide may be purified, isolated or enriched using a variety of methods. In some methods, the protein or polypeptide may be secreted into the culture medium via a native signal peptide or a heterologous signal peptide operably linked thereto. In some methods, the protein or polypeptide may be linked to a heterologous polypeptide which facilitates its isolation,  
20 purification, or enrichment such as a nickel binding polypeptide. The protein or polypeptide may also be obtained by gel electrophoresis, ion exchange chromatography, size chromatography, hplc, salt precipitation, immunoprecipitation, a combination of any of the preceding methods, or any of the isolation, purification, or enrichment techniques familiar to those skilled in the art.

The protein encoded by the nucleic acid insert may also be purified using standard  
25 immunochromatography techniques using immunoaffinity chromatography with antibodies directed against the encoded protein or polypeptide as described in more detail below. If antibody production is not possible, the nucleic acid insert encoding the desired protein or polypeptide may be incorporated into expression vectors designed for use in purification schemes employing chimeric polypeptides. In such strategies, the coding sequence of the nucleic acid insert is ligated in frame with the gene encoding the  
30 other half of the chimera. The other half of the chimera may be  $\beta$ -globin or a nickel binding polypeptide. A chromatography matrix having antibody to  $\beta$ -globin or nickel attached thereto is then used to purify the chimeric protein. Protease cleavage sites may be engineered between the  $\beta$ -globin gene or the nickel binding polypeptide and the extended cDNA or portion thereof. Thus, the two polypeptides of the chimera may be separated from one another by protease digestion.

35 One useful expression vector for generating  $\beta$ -globin chimerics is pSG5 (Stratagene), which encodes rabbit  $\beta$ -globin. Intron II of the rabbit  $\beta$ -globin gene facilitates splicing of the expressed transcript, and the polyadenylation signal incorporated into the construct increases the level of

expression. These techniques as described are well known to those skilled in the art of molecular biology. Standard methods are published in methods texts such as Davis *et al.*, (*Basic Methods in Molecular Biology*, L.G. Davis, M.D. Digner, and J.F. Battey, ed., Elsevier Press, NY, 1986) and many of the methods are available from Stratagene, Life Technologies, Inc., or Promega. Polypeptide may additionally be produced from the construct using *in vitro* translation systems such as the *In vitro* Express™ Translation Kit (Stratagene).

Following expression and purification of the proteins or polypeptides encoded by the nucleic acid inserts, the purified proteins may be tested for the ability to bind to the surface of various cell types as described in Example 23 below. It will be appreciated that a plurality of proteins expressed from these nucleic acid inserts may be included in a panel of proteins to be simultaneously evaluated for the activities specifically described below, as well as other biological roles for which assays for determining activity are available.

### EXAMPLE 23

#### Analysis of Secreted Proteins to Determine Whether they Bind to the Cell Surface

The EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of EST-related nucleic acids, fragments of positional segments of EST-related nucleic acids, nucleic acids encoding the EST-related polypeptides, nucleic acids encoding fragments of the EST-related polypeptides, nucleic acids encoding positional segments of EST-related polypeptides, or nucleic acids encoding fragments of positional segments of EST-related polypeptides are cloned into expression vectors such as those described in Example 22. The encoded proteins or polypeptides are purified, isolated, or enriched as described above. Following purification, isolation, or enrichment, the proteins or polypeptides are labeled using techniques known to those skilled in the art. The labeled proteins or polypeptides are incubated with cells or cell lines derived from a variety of organs or tissues to allow the proteins to bind to any receptor present on the cell surface. Following the incubation, the cells are washed to remove non-specifically bound proteins or polypeptides. The specifically bound labeled proteins or polypeptides are detected by autoradiography. Alternatively, unlabeled proteins or polypeptides may be incubated with the cells and detected with antibodies having a detectable label, such as a fluorescent molecule, attached thereto.

Specificity of cell surface binding may be analyzed by conducting a competition analysis in which various amounts of unlabeled protein or polypeptide are incubated along with the labeled protein or polypeptide. The amount of labeled protein or polypeptide bound to the cell surface decreases as the amount of competitive unlabeled protein or polypeptide increases. As a control, various amounts of an unlabeled protein or polypeptide unrelated to the labeled protein or polypeptide is included in some binding reactions. The amount of labeled protein or polypeptide bound to the cell surface does not decrease in binding reactions containing increasing amounts of unrelated unlabeled protein, indicating that the protein or polypeptide encoded by the nucleic acid binds specifically to the cell surface.

As discussed above, human proteins have been shown to have a number of important physiological effects and, consequently, represent a valuable therapeutic resource. The human proteins or polypeptides made as described above may be evaluated to determine their physiological activities as described below.

5

#### EXAMPLE 24

##### Assaying the Expressed Proteins or Polypeptides for Cytokine,

##### Cell Proliferation or Cell Differentiation Activity

As discussed above, some human proteins act as cytokines or may affect cellular proliferation or differentiation. Many protein factors discovered to date, including all known cytokines, have exhibited activity in one or more factor dependent cell proliferation assays, and hence the assays serve as a convenient confirmation of cytokine activity. The activity of a protein or polypeptide of the present invention is evidenced by any one of a number of routine factor dependent cell proliferation assays for cell lines including, without limitation, 32D, DA2, DA1G, T10, B9, B9/11, BaF3, MC9/G, M<sup>+</sup> (preB  
10 M<sup>+</sup>), 2E8, RB5, DA1, 123, T1165, HT2, CTLL2, TF-1, Mo7c and CMK. The proteins or polypeptides prepared as described above may be evaluated for their ability to regulate T cell or thymocyte proliferation in assays such as those described above or in the following references: *Current Protocols in Immunology*, Ed. by J.E. Coligan *et al.*, Greene Publishing Associates and Wiley-Interscience; Takai *et al. J. Immunol.* 137:3494-3500, 1986., Bertagnolli *et al. J. Immunol.* 145:1706-1712, 1990.,  
15 Bertagnolli *et al., Cellular Immunology* 133:327-341, 1991. Bertagnolli, *et al. J. Immunol.* 149:3778-3783, 1992; Bowman *et al., J. Immunol.* 152:1756-1761, 1994.

In addition, numerous assays for cytokine production and/or the proliferation of spleen cells, lymph node cells and thymocytes are known. These include the techniques disclosed in *Current Protocols in Immunology*. J.E. Coligan *et al.* Eds., 1:3.12.1-3.12.14, John Wiley and Sons, Toronto.  
25 1994; and Schreiber, R.D. In *Current Protocols in Immunology.*, *supra* 1 : 6.8.1-6.8.8.

The proteins or polypeptides prepared as described above may also be assayed for the ability to regulate the proliferation and differentiation of hematopoietic or lymphopoietic cells. Many assays for such activity are familiar to those skilled in the art, including the assays in the following references: Bottomly *et al.*, In *Current Protocols in Immunology.*, *supra* 1 : 6.3.1-6.3.12.; deVries *et al., J. Exp.*  
30 *Med.* 173:1205-1211, 1991; Moreau *et al., Nature* 36:690-692, 1988; Greenberger *et al., Proc. Natl. Acad. Sci. U.S.A.* 80:2931-2938, 1983; Nordan, R., In *Current Protocols in Immunology.*, *supra* 1 : 6.6.1-6.6.5; Smith *et al., Proc. Natl. Acad. Sci. U.S.A.* 83:1857-1861, 1986; Bennett *et al* in *Current Protocols in Immunology supra* 1 : 6.15.1; Ciarletta *et al* In *Current Protocols in Immunology. supra* 1 : 6.13.1.

35 The proteins or polypeptides prepared as described above may also be assayed for their ability to regulate T-cell responses to antigens. Many assays for such activity are familiar to those skilled in the art, including the assays described in the following references: Chapter 3 (*In vitro* Assays for Mouse

Lymphocyte Function), Chapter 6 (Cytokines and Their Cellular Receptors) and Chapter 7, (Immunologic Studies in Humans) in *Current Protocols in Immunology supra*; Weinberger *et al.*, *Proc. Natl. Acad. Sci. USA* 77:6091-6095, 1980; Weinberger *et al.*, *Eur. J. Immunol.* 11:405-411, 1981; Takai *et al.*, *J. Immunol.* 137:3494-3500, 1986; Takai *et al.*, *J. Immunol.* 140:508-512, 1988.

5        Those proteins or polypeptides which exhibit cytokine, cell proliferation, or cell differentiation activity may then be formulated as pharmaceuticals and used to treat clinical conditions in which induction of cell proliferation or differentiation is beneficial. Alternatively, as described in more detail below, nucleic acids encoding these proteins or polypeptides or nucleic acids regulating the expression of these proteins or polypeptides may be introduced into appropriate host cells to increase or decrease the  
10 expression of the proteins or polypeptides as desired.

### EXAMPLE 25

#### Assaying the Expressed Proteins or Polypeptides for Activity as Immune System Regulators

15        The proteins or polypeptides prepared as described above may also be evaluated for their effects as immune regulators. For example, the proteins or polypeptides may be evaluated for their activity to influence thymocyte or splenocyte cytotoxicity. Numerous assays for such activity are familiar to those skilled in the art including the assays described in the following references: Chapter 3 (*In vitro* Assays for Mouse Lymphocyte Function 3.1-3.19) and Chapter 7 (Immunologic studies in Humans) in *Current  
20 Protocols in Immunology*, J.E. Coligan *et al.* Eds, Greene Publishing Associates and Wiley-Interscience; Herrmann *et al.*, *Proc. Natl. Acad. Sci. USA* 78:2488-2492, 1981; Herrmann *et al.*, *J. Immunol.* 128:1968-1974, 1982; Handa *et al.*, *J. Immunol.* 135:1564-1572, 1985; Takai *et al.*, *J. Immunol.* 137:3494-3500, 1986; Takai *et al.*, *J. Immunol.* 140:508-512, 1988; Bowman *et al.*, *J. Virology* 61:1992-1998; Bertagnolli *et al. Cell. Immunol.* 133:327-341, 1991; Brown *et al.*, *J. Immunol.* 153:3079-3092,  
25 1994.

      The proteins or polypeptides prepared as described above may also be evaluated for their effects on T-cell dependent immunoglobulin responses and isotype switching. Numerous assays for such activity are familiar to those skilled in the art, including the assays disclosed in the following references: Maliszewski, *J. Immunol.* 144:3028-3033, 1990; Mond *et al.* in *Current Protocols in Immunology*, 1 :  
30 3.8.1-3.8.16, *supra*.

      The proteins or polypeptides prepared as described above may also be evaluated for their effect on immune effector cells, including their effect on Th1 cells and cytotoxic lymphocytes. Numerous assays for such activity are familiar to those skilled in the art, including the assays disclosed in the following references: Chapter 3 (*In vitro* Assays for Mouse Lymphocyte Function 3.1-3.19) and Chapter  
35 7 (Immunologic Studies in Humans) in *Current Protocols in Immunology, supra*; Takai *et al.*, *J. Immunol.* 137:3494-3500, 1986; Takai *et al.*, *J. Immunol.* 140:508-512, 1988; Bertagnolli *et al.*, *J. Immunol.* 149:3778-3783, 1992.

The proteins or polypeptides prepared as described above may also be evaluated for their effect on dendritic cell mediated activation of naive T-cells. Numerous assays for such activity are familiar to those skilled in the art, including the assays disclosed in the following references: Guery *et al.*, *J. Immunol.* 134:536-544, 1995; Inaba *et al.*, *J. Exp. Med.* 173:549-559, 1991; Macatonia *et al.*, *J. Immunol.* 154:5071-5079, 1995; Porgador *et al.* *J. Exp. Med.* 182:255-260, 1995; Nair *et al.*, *J. Virol.* 67:4062-4069, 1993; Huang *et al.*, *Science* 264:961-965, 1994; Macatonia *et al.* *J. Exp. Med.* 169:1255-1264, 1989; Bhardwaj *et al.*, *Journal of Clinical Investigation* 94:797-807, 1994; and Inaba *et al.*, *J. Exp. Med.* 172:631-640, 1990.

The proteins or polypeptides prepared as described above may also be evaluated for their influence on the lifetime of lymphocytes. Numerous assays for such activity are familiar to those skilled in the art, including the assays disclosed in the following references: Darzynkiewicz *et al.*, *Cytometry* 13:795-808, 1992; Gorczyca *et al.*, *Leukemia* 7:659-670, 1993; Gorczyca *et al.*, *Cancer Res.* 53:1945-1951, 1993; Itoh *et al.*, *Cell* 66:233-243, 1991; Zacharchuk, *J. Immunol.* 145:4037-4045, 1990; Zamai *et al.*, *Cytometry* 14:891-897, 1993; Gorczyca *et al.*, *Int. J. Oncol.* 1:639-648, 1992.

The proteins or polypeptides prepared as described above may also be evaluated for their influence on early steps of T-cell commitment and development. Numerous assays for such activity are familiar to those skilled in the art, including without limitation the assays disclosed in the following references: Antica *et al.*, *Blood* 84:111-117, 1994; Fine *et al.*, *Cell. Immunol.* 155:111-122, 1994; Galy *et al.*, *Blood* 85:2770-2778, 1995; Toki *et al.*, *Proc. Nat. Acad. Sci. USA* 88:7548-7551, 1991.

Those proteins or polypeptides which exhibit activity as immune system regulators activity may then be formulated as pharmaceuticals and used to treat clinical conditions in which regulation of immune activity is beneficial. For example, the protein or polypeptide may be useful in the treatment of various immune deficiencies and disorders (including severe combined immunodeficiency), e.g., in regulating (up or down) growth and proliferation of T and/or B lymphocytes, as well as effecting the cytolytic activity of NK cells and other cell populations. These immune deficiencies may be genetic or be caused by viral (e.g., HIV) as well as bacterial or fungal infections, or may result from autoimmune disorders. More specifically, infectious diseases caused by viral, bacterial, fungal or other infection may be treatable using the protein or polypeptide including infections by HIV, hepatitis viruses, herpesviruses, mycobacteria, *Leishmania* spp., *plamodium*. and various fungal infections such as candidiasis. Of course, in this regard, a protein or polypeptide may also be useful where a boost to the immune system generally may be desirable, *i.e.*, in the treatment of cancer.

Alternatively, the proteins or polypeptides prepared as described above may be used in treatment of autoimmune disorders including, for example, connective tissue disease, multiple sclerosis, systemic lupus erythematosus, rheumatoid arthritis, autoimmune pulmonary inflammation, Guillain-Barre syndrome, autoimmune thyroiditis, insulin dependent diabetes mellitis, myasthenia gravis, graft-versus-host disease and autoimmune inflammatory eye disease. Such a protein or polypeptide may also be useful in the treatment of allergic reactions and conditions, such as asthma (particularly allergic

asthma) or other respiratory problems. Other conditions, in which immune suppression is desired (including, for example, organ transplantation), may also be treatable using the protein or polypeptide.

Using the proteins or polypeptides of the invention it may also be possible to regulate immune responses either up or down. Down regulation may involve inhibiting or blocking an immune response  
5 already in progress or may involve preventing the induction of an immune response. The functions of activated T-cells may be inhibited by suppressing T cell responses or by inducing specific tolerance in T cells, or both. Immunosuppression of T cell responses is generally an active non-antigen-specific process which requires continuous exposure of the T cells to the suppressive agent. Tolerance, which involves inducing non-responsiveness or anergy in T cells, is distinguishable from immunosuppression  
10 in that it is generally antigen-specific and persists after the end of exposure to the tolerizing agent. Operationally, tolerance can be demonstrated by the lack of a T cell response upon reexposure to specific antigen in the absence of the tolerizing agent.

Down regulating or preventing one or more antigen functions (including without limitation B lymphocyte antigen functions, such as, for example, B7 costimulation), e.g., preventing high level  
15 lymphokine synthesis by activated T cells, will be useful in situations of tissue, skin and organ transplantation and in graft-versus-host disease (GVHD). For example, blockage of T cell function should result in reduced tissue destruction in tissue transplantation. Typically, in tissue transplants, rejection of the transplant is initiated through its recognition as foreign by T cells, followed by an immune reaction that destroys the transplant. The administration of a molecule which inhibits or blocks  
20 interaction of a B7 lymphocyte antigen with its natural ligand(s) on immune cells (such as a soluble, monomeric form of a peptide having B7-2 activity alone or in conjunction with a monomeric form of a peptide having an activity of another B lymphocyte antigen (e.g., B7-1, B7-3) or blocking antibody), prior to transplantation, can lead to the binding of the molecule to the natural ligand(s) on the immune cells without transmitting the corresponding costimulatory signal. Blocking B lymphocyte antigen  
25 function in this matter prevents cytokine synthesis by immune cells, such as T cells, and thus acts as an immunosuppressant. Moreover, the lack of costimulation may also be sufficient to anergize the T cells, thereby inducing tolerance in a subject. Induction of long-term tolerance by B lymphocyte antigen-blocking reagents may avoid the necessity of repeated administration of these blocking reagents. To achieve sufficient immunosuppression or tolerance in a subject, it may also be necessary to block the  
30 function of a combination of B lymphocyte antigens.

The efficacy of particular blocking reagents in preventing organ transplant rejection or GVHD can be assessed using animal models that are predictive of efficacy in humans. Examples of appropriate systems which can be used include allogeneic cardiac grafts in rats and xenogeneic pancreatic islet cell grafts in mice, both of which have been used to examine the immunosuppressive effects of CTLA4Ig  
35 fusion proteins *in vivo* as described in Lenschow *et al.*, *Science* 257:789-792 (1992) and Turka *et al.*, *Proc. Natl. Acad. Sci USA*, 89:11102-11105 (1992). In addition, murine models of GVHD (see Paul *et al.*,

*Fundamental Immunology*, Raven Press, New York, 1989, pp. 846-847) can be used to determine the effect of blocking B lymphocyte antigen function *in vivo* on the development of that disease.

Blocking antigen function may also be therapeutically useful for treating autoimmune diseases. Many autoimmune disorders are the result of inappropriate activation of T cells that are reactive against self tissue and which promote the production of cytokines and autoantibodies involved in the pathology of the diseases. Preventing the activation of autoreactive T cells may reduce or eliminate disease symptoms. Administration of reagents which block costimulation of T cells by disrupting receptor/ligand interactions of B lymphocyte antigens can be used to inhibit T cell activation and prevent production of autoantibodies or T cell-derived cytokines which potentially involved in the disease process. Additionally, blocking reagents may induce antigen-specific tolerance of autoreactive T cells which could lead to long-term relief from the disease. The efficacy of blocking reagents in preventing or alleviating autoimmune disorders can be determined using a number of well-characterized animal models of human autoimmune diseases. Examples include murine experimental autoimmune encephalitis, systemic lupus erythmatosis in MRL/pr/pr mice or NZB hybrid mice, murine autoimmune collagen arthritis, diabetes mellitus in OD mice and BB rats, and murine experimental myasthenia gravis (see Paul ed., *Fundamental Immunology*, Raven Press, New York, 1989, pp. 840-856).

Upregulation of an antigen function (preferably a B lymphocyte antigen function), as a means of up regulating immune responses, may also be useful in therapy. Upregulation of immune responses may involve either enhancing an existing immune response or eliciting an initial immune response as shown by the following examples. For instance, enhancing an immune response through stimulating B lymphocyte antigen function may be useful in cases of viral infection. In addition, systemic viral diseases such as influenza, the common cold, and encephalitis might be alleviated by the administration of stimulatory form of B lymphocyte antigens systemically.

Alternatively, antiviral immune responses may be enhanced in an infected patient by removing T cells from the patient, costimulating the T cells *in vitro* with viral antigen-pulsed APCs either expressing the proteins or polypeptides described above or together with a stimulatory form of the protein or polypeptide and reintroducing the *in vitro* primed T cells into the patient. The infected cells would now be capable of delivering a costimulatory signal to T cells *in vivo*, thereby activating the T cells.

In another application, upregulation or enhancement of antigen function (preferably B lymphocyte antigen function) may be useful in the induction of tumor immunity. Tumor cells (e.g., sarcoma, melanoma, lymphoma, leukemia, neuroblastoma, carcinoma) transfected with one of the above-described nucleic acids encoding a protein or polypeptide can be administered to a subject to overcome tumor-specific tolerance in the subject. If desired, the tumor cell can be transfected to express a combination of peptides. For example, tumor cells obtained from a patient can be transfected *ex vivo* with an expression vector directing the expression of a peptide having B7-2-like activity alone, or in conjunction with a peptide having B7-1-like activity and/or B7-3-like activity. The transfected tumor

cells are returned to the patient to result in expression of the peptides on the surface of the transfected cell. Alternatively, gene therapy techniques can be used to target a tumor cell for transfection *in vivo*.

The presence of the protein or polypeptide encoded by the nucleic acids described above having the activity of a B lymphocyte antigen(s) on the surface of the tumor cell provides the necessary

5 costimulation signal to T cells to induce a T cell mediated immune response against the transfected tumor cells. In addition, tumor cells which lack or which fail to reexpress sufficient amounts of MHC class I or MHC class II molecules can be transfected with nucleic acids encoding all or a portion of (e.g., a cytoplasmic-domain truncated portion) of an MHC class I  $\alpha$  chain and  $\beta_2$  microglobulin or an MHC class II  $\alpha$  chain and an MHC class II  $\beta$  chain to thereby express MHC class I or MHC class II proteins

10 on the cell surface, respectively. Expression of the appropriate MHC class I or class II molecules in conjunction with a peptide having the activity of a B lymphocyte antigen (e.g., B7-1, B7-2, B7-3) induces a T cell mediated immune response against the transfected tumor cell. Optionally, a nucleic acid encoding an antisense construct which blocks expression of an MHC class II associated protein, such as the invariant chain, can also be cotransfected with a DNA encoding a protein or polypeptide having the

15 activity of a B lymphocyte antigen to promote presentation of tumor associated antigens and induce tumor specific immunity. Thus, the induction of a T cell mediated immune response in a human subject may be sufficient to overcome tumor-specific tolerance in the subject. Alternatively, as described in more detail below, nucleic acids encoding these immune system regulator proteins or polypeptides or nucleic acids regulating the expression of such proteins or polypeptides may be introduced into

20 appropriate host cells to increase or decrease the expression of the proteins as desired.

### EXAMPLE 26

#### Assaying the Expressed Proteins or Polypeptides for Hematopoiesis Regulating Activity

25 The proteins or polypeptides encoded by the nucleic acids described above may also be evaluated for their hematopoiesis regulating activity. For example, the effect of the proteins or polypeptides on embryonic stem cell differentiation may be evaluated. Numerous assays for such activity are familiar to those skilled in the art, including the assays disclosed in the following references: Johansson *et al.*, *Cell. Biol.* 15:141-151, 1995; Keller *et al.*, *Mol. Cell. Biol.* 13:473-486, 1993;

30 McClanahan *et al.*, *Blood* 81:2903-2915, 1993.

The proteins or polypeptides encoded by the nucleic acids described above may also be evaluated for their influence on the lifetime of stem cells and stem cell differentiation. Numerous assays for such activity are familiar to those skilled in the art, including the assays disclosed in the following references: Freshney, M.G. Methylcellulose Colony Forming Assays, in Culture of Hematopoietic Cells.

35 R.I. Freshney, *et al.* Eds. pp. 265-268, Wiley-Liss, Inc., New York, NY. 1994; Hirayama *et al.*, *Proc. Natl. Acad. Sci. USA* 89:5907-5911, 1992; McNiece, I.K. and Briddell, R.A. Primitive Hematopoietic Colony Forming Cells with High Proliferative Potential, in Culture of Hematopoietic Cells. supra;



Neben *et al.*, *Experimental Hematology* 22:353-359, 1994; Ploemacher, R.E. Cobblestone Area Forming Cell Assay, In Culture of Hematopoietic Cells, *supra*; Spooncer, E., Dexter, M. and Allen, T. Long Term Bone Marrow Cultures in the Presence of Stromal Cells, in Culture of Hematopoietic Cells *supra*; and Sutherland, H.J. Long Term Culture Initiating Cell Assay, in Culture of Hematopoietic Cells, *supra*.

5        Those proteins or polypeptides which exhibit hematopoiesis regulatory activity may then be formulated as pharmaceuticals and used to treat clinical conditions in which regulation of hematopoiesis is beneficial. For example, a protein or polypeptide of the present invention may be useful in regulation of hematopoiesis and, consequently, in the treatment of myeloid or lymphoid cell deficiencies. Even marginal biological activity in support of colony forming cells or of factor-dependent cell lines indicates  
10 involvement in regulating hematopoiesis, e.g. in supporting the growth and proliferation of erythroid progenitor cells alone or in combination with other cytokines, thereby indicating utility, for example, in treating various anemias or for use in conjunction with irradiation/chemotherapy to stimulate the production of erythroid precursors and/or erythroid cells; in supporting the growth and proliferation of myeloid cells such as granulocytes and monocytes/macrophages (*i.e.*, traditional CSF activity) useful, for  
15 example, in conjunction with chemotherapy to prevent or treat consequent myelo-suppression; in supporting the growth and proliferation of megakaryocytes and consequently of platelets thereby allowing prevention or treatment of various platelet disorders such as thrombocytopenia, and generally for use in place of or complimentary to platelet transfusions; and/or in supporting the growth and proliferation of hematopoietic stem cells which are capable of maturing to any and all of the above-  
20 mentioned hematopoietic cells and therefore find therapeutic utility in various stem cell disorders (such as those usually treated with transplantation, including, without limitation, aplastic anemia and paroxysmal nocturnal hemoglobinuria), as well as in repopulating the stem cell compartment post irradiation/chemotherapy, either in-vivo or ex-vivo (*i.e.*, in conjunction with bone marrow transplantation or with peripheral progenitor cell transplantation (homologous or heterologous)) as  
25 normal cells or genetically manipulated for gene therapy. Alternatively, as described in more detail below, nucleic acids encoding these proteins or polypeptides or nucleic acids regulating the expression of these proteins or polypeptides may be introduced into appropriate host cells to increase or decrease the expression of the proteins as desired.

30

#### EXAMPLE 27

##### Assaying the Expressed Proteins or Polypeptides for Regulation of Tissue Growth

The proteins or polypeptides encoded by the nucleic acids described above may also be evaluated for their effect on tissue growth. Numerous assays for such activity are familiar to those  
35 skilled in the art, including the assays disclosed in International Patent Publication No. WO95/16035, International Patent Publication No. WO95/05846 and International Patent Publication No. WO91/07491.

Assays for wound healing activity include, without limitation, those described in: Winter, *Epidermal Wound Healing*, pps. 71-112 (Maibach, H1 and Rovee, DT, eds.), Year Book Medical Publishers, Inc., Chicago, as modified by Eaglstein and Mertz, J. Invest. Dermatol 71:382-84 (1978).

Those proteins or polypeptides which are involved in the regulation of tissue growth may then  
5 be formulated as pharmaceuticals and used to treat clinical conditions in which regulation of tissue growth is beneficial. For example, a protein or polypeptide may have utility in compositions used for bone, cartilage, tendon, ligament and/or nerve tissue growth or regeneration, as well as for wound healing and tissue repair and replacement, and in the treatment of burns, incisions and ulcers.

A protein or polypeptide encoded by the nucleic acids described above which induces cartilage  
10 and/or bone growth in circumstances where bone is not normally formed, has application in the healing of bone fractures and cartilage damage or defects in humans and other animals. Such a preparation employing a protein or polypeptide of the invention may have prophylactic use in closed as well as open fracture reduction and also in the improved fixation of artificial joints. *De novo* bone synthesis induced by an osteogenic agent contributes to the repair of congenital, trauma induced, or oncologic resection  
15 induced craniofacial defects, and also is useful in cosmetic plastic surgery.

A protein or polypeptide of this invention may also be used in the treatment of periodontal disease, and in other tooth repair processes. Such agents may provide an environment to attract bone-forming cells, stimulate growth of bone-forming cells or induce differentiation of progenitors of bone-forming cells. A protein of the invention may also be useful in the treatment of osteoporosis or  
20 osteoarthritis, such as through stimulation of bone and/or cartilage repair or by blocking inflammation or processes of tissue destruction (collagenase activity, osteoclast activity, etc.) mediated by inflammatory processes.

Another category of tissue regeneration activity that may be attributable to the proteins or polypeptides encoded by the nucleic acids described above is tendon/ligament formation. A protein or  
25 polypeptide encoded by the nucleic acids described above, which induces tendon/ligament-like tissue or other tissue formation in circumstances where such tissue is not normally formed, has application in the healing of tendon or ligament tears, deformities and other tendon or ligament defects in humans and other animals. Such a preparation employing a tendon/ligament-like tissue inducing protein may have prophylactic use in preventing damage to tendon or ligament tissue, as well as use in the improved  
30 fixation of tendon or ligament to bone or other tissues, and in repairing defects to tendon or ligament tissue. *De novo* tendon/ligament-like tissue formation induced by a protein or polypeptide of the present invention contributes to the repair of tendon or ligaments defects of congenital, traumatic or other origin and is also useful in cosmetic plastic surgery for attachment or repair of tendons or ligaments. The proteins or polypeptides of the present invention may provide an environment to attract tendon- or  
35 ligament-forming cells, stimulate growth of tendon- or ligament-forming cells, induce differentiation of progenitors of tendon- or ligament-forming cells, or induce growth of tendon/ligament cells or progenitors *ex vivo* for return *in vivo* to effect tissue repair. The proteins or polypeptides of the

invention may also be useful in the treatment of tendinitis, carpal tunnel syndrome and other tendon or ligament defects. The therapeutic compositions may also include an appropriate matrix and/or sequestering agent as a carrier as is well known in the art.

The proteins or polypeptides of the present invention may also be useful for proliferation of  
5 neural cells and for regeneration of nerve and brain tissue, *i.e.*, for the treatment of central and peripheral nervous system diseases and neuropathies, as well as mechanical and traumatic disorders, which involve degeneration, death or trauma to neural cells or nerve tissue. More specifically, a protein or polypeptide may be used in the treatment of diseases of the peripheral nervous system, such as peripheral nerve injuries, peripheral neuropathy and localized neuropathies, and central nervous system diseases, such as  
10 Alzheimer's, Parkinson's disease, Huntington's disease, amyotrophic lateral sclerosis, and Shy-Drager syndrome. Further conditions which may be treated in accordance with the present invention include mechanical and traumatic disorders, such as spinal cord disorders, head trauma and cerebrovascular diseases such as stroke. Peripheral neuropathies resulting from chemotherapy or other medical therapies may also be treatable using a protein or polypeptide of the invention.

15 Proteins or polypeptides of the invention may also be useful to promote better or faster closure of non-healing wounds, including without limitation pressure ulcers, ulcers associated with vascular insufficiency, surgical and traumatic wounds, and the like.

It is expected that a protein or polypeptide of the present invention may also exhibit activity for generation or regeneration of other tissues, such as organs (including, for example, pancreas, liver,  
20 intestine, kidney, skin, endothelium) muscle (smooth, skeletal or cardiac) and vascular (including vascular endothelium) tissue, or for promoting the growth of cells comprising such tissues. Part of the desired effects may be by inhibition or modulation of fibrotic scarring to allow normal tissue to generate. A protein or polypeptide of the invention may also exhibit angiogenic activity.

A protein or polypeptide of the present invention may also be useful for gut protection or  
25 regeneration and treatment of lung or liver fibrosis, reperfusion injury in various tissues, and conditions resulting from systemic cytokine damage.

A protein or polypeptide of the present invention may also be useful for promoting or inhibiting differentiation of tissues described above from precursor tissues or cells; or for inhibiting the growth of tissues described above.

30 Alternatively, as described in more detail below, nucleic acids encoding tissue growth regulating activity proteins or polypeptides or nucleic acids regulating the expression of such proteins or polypeptides may be introduced into appropriate host cells to increase or decrease the expression of the proteins as desired.

35

#### EXAMPLE 28

##### Assaying the Expressed Proteins or Polypeptides for Regulation of Reproductive Hormones

The proteins or polypeptides of the present invention may also be evaluated for their ability to regulate reproductive hormones, such as follicle stimulating hormone. Numerous assays for such activity are familiar to those skilled in the art, including the assays disclosed in the following references: Vale *et al.*, *Endocrinol.* 91:562-572, 1972; Ling *et al.*, *Nature* 321:779-782, 1986; Vale *et al.*, *Nature* 321:776-779, 1986; Mason *et al.*, *Nature* 318:659-663, 1985; Forage *et al.*, *Proc. Natl. Acad. Sci. USA* 83:3091-3095, 1986. Chapter 6.12 in *Current Protocols in Immunology*, J.E. Coligan *et al.* Eds. Greene Publishing Associates and Wiley-Interscience; Taub *et al.* *J. Clin. Invest.* 95:1370-1376, 1995; Lind *et al.* *APMIS* 103:140-146, 1995; Muller *et al.* *Eur. J. Immunol.* 25:1744-1748; Gruber *et al.* *J. Immunol.* 152:5860-5867, 1994; Johnston *et al.*, *J Immunol.* 153:1762-1768, 1994.

Those proteins or polypeptides which exhibit activity as reproductive hormones or regulators of cell movement may then be formulated as pharmaceuticals and used to treat clinical conditions in which regulation of reproductive hormones are beneficial. For example, a protein or polypeptide may exhibit activin- or inhibin-related activities. Inhibins are characterized by their ability to inhibit the release of follicle stimulating hormone (FSH), while activins are characterized by their ability to stimulate the release of FSH. Thus, a protein or polypeptide of the present invention, alone or in heterodimers with a member of the inhibin family, may be useful as a contraceptive based on the ability of inhibins to decrease fertility in female mammals and decrease spermatogenesis in male mammals. Administration of sufficient amounts of other inhibins can induce infertility in these mammals. Alternatively, the protein or polypeptide of the invention, as a homodimer or as a heterodimer with other protein subunits of the inhibin-B group, may be useful as a fertility inducing therapeutic, based upon the ability of activin molecules in stimulating FSH release from cells of the anterior pituitary. See, for example, United States Patent 4,798,885. A protein or polypeptide of the invention may also be useful for advancement of the onset of fertility in sexually immature mammals, so as to increase the lifetime reproductive performance of domestic animals such as cows, sheep and pigs.

Alternatively, as described in more detail below, nucleic acids encoding reproductive hormone regulating activity proteins or polypeptides or nucleic acids regulating the expression of such proteins or polypeptides may be introduced into appropriate host cells to increase or decrease the expression of the proteins or polypeptides as desired.

## EXAMPLE 29

### Assaying the Expressed Proteins or Polypeptides For Chemotactic/Chemokinetic Activity

The proteins or polypeptides of the present invention may also be evaluated for chemotactic/chemokinetic activity. For example, a protein or polypeptide of the present invention may have chemotactic or chemokinetic activity (e.g., act as a chemokine) for mammalian cells, including, for example, monocytes, fibroblasts, neutrophils, T-cells, mast cells, eosinophils, epithelial and/or endothelial cells. Chemotactic and chemokinetic proteins or polypeptides can be used to mobilize or attract a desired cell population to a desired site of action. Chemotactic or chemokinetic proteins or

polypeptides provide particular advantages in treatment of wounds and other trauma to tissues, as well as in treatment of localized infections. For example, attraction of lymphocytes, monocytes or neutrophils to tumors or sites of infection may result in improved immune responses against the tumor or infecting agent.

5 A protein or polypeptide has chemotactic activity for a particular cell population if it can stimulate, directly or indirectly, the directed orientation or movement of such cell population. Preferably, the protein or polypeptide has the ability to directly stimulate directed movement of cells. Whether a particular protein or polypeptide has chemotactic activity for a population of cells can be readily determined by employing such protein or polypeptide in any known assay for cell chemotaxis.

10 The activity of a protein or polypeptide of the invention may, among other means, be measured by the following methods:

Assays for chemotactic activity (which will identify proteins or polypeptides that induce or prevent chemotaxis) consist of assays that measure the ability of a protein or polypeptide to induce the migration of cells across a membrane as well as the ability of a protein or polypeptide to induce the  
15 adhesion of one cell population to another cell population. Suitable assays for movement and adhesion include, without limitation, those described in: *Current Protocols in Immunology*, Ed by J.E. Coligan, A.M. Kruisbeek, D.H. Margulies, E.M. Shevach, W. Strober, Pub. Greene Publishing Associates and Wiley-Interscience, Chapter 6.12: 6.12.1-6.12.28; Taub *et al. J. Clin. Invest.* 95:1370-1376, 1995; Lind *et al. APMIS* 103:140-146, 1995; Mueller *et al., Eur. J. Immunol.* 25:1744-1748; Gruber *et al. J. Immunol.*  
20 152:5860-5867, 1994; Johnston *et al. J. Immunol.*, 153:1762-1768, 1994.

### EXAMPLE 30

#### Assaying the Expressed Proteins or Polypeptides for Regulation of Blood Clotting

The proteins or polypeptides of the present invention may also be evaluated for their effects on  
25 blood clotting. Numerous assays for such activity are familiar to those skilled in the art, including the assays disclosed in the following references: Linet *et al., J. Clin. Pharmacol.* 26:131-140, 1986; Burdick *et al., Thrombosis Res.* 45:413-419, 1987; Humphrey *et al., Fibrinolysis* 5:71-79 (1991); Schaub, *Prostaglandins* 35:467-474, 1988.

Those proteins or polypeptides which are involved in the regulation of blood clotting may then  
30 be formulated as pharmaceuticals and used to treat clinical conditions in which regulation of blood clotting is beneficial. For example, a protein or polypeptide of the invention may also exhibit hemostatic or thrombolytic activity. As a result, such a protein or polypeptide is expected to be useful in treatment of various coagulations disorders (including hereditary disorders, such as hemophilias) or to enhance coagulation and other hemostatic events in treating wounds resulting from trauma, surgery or other  
35 causes. A protein or polypeptide of the invention may also be useful for dissolving or inhibiting formation of thromboses and for treatment and prevention of conditions resulting therefrom (such as infarction of cardiac and central nervous system vessels (e.g., stroke)). Alternatively, as described in

more detail below, nucleic acids encoding blood clotting activity proteins or polypeptides or nucleic acids regulating the expression of such proteins or polypeptides may be introduced into appropriate host cells to increase or decrease the expression of the proteins or polypeptides as desired.

5

### EXAMPLE 31

#### Assaying the Expressed Proteins or Polypeptides for Involvement in Receptor/Ligand Interactions

The proteins or polypeptides of the present invention may also be evaluated for their involvement in receptor/ligand interactions. Numerous assays for such involvement are familiar to those skilled in the art, including the assays disclosed in the following references: Chapter 7. 7.28.1-7.28.22) in *Current Protocols in Immunology*, J.E. Coligan *et al.* Eds. Greene Publishing Associates and Wiley-Interscience; Takai *et al.*, *Proc. Natl. Acad. Sci. USA* 84:6864-6868, 1987; Bierer *et al.*, *J. Exp. Med.* 168:1145-1156, 1988; Rosenstein *et al.*, *J. Exp. Med.* 169:149-160, 1989; Stoltenborg *et al.*, *J. Immunol. Methods* 175:59-68, 1994; Stitt *et al.*, *Cell* 80:661-670, 1995; Gyuris *et al.*, *Cell* 75:791-803, 1993.

For example, the proteins or polypeptides of the present invention may also demonstrate activity as receptors, receptor ligands or inhibitors or agonists of receptor/ligand interactions. Examples of such receptors and ligands include, without limitation, cytokine receptors and their ligands, receptor kinases and their ligands, receptor phosphatases and their ligands, receptors involved in cell-cell interactions and their ligands (including without limitation, cellular adhesion molecules (such as selectins, integrins and their ligands) and receptor/ligand pairs involved in antigen presentation, antigen recognition and development of cellular and humoral immune responses). Receptors and ligands are also useful for screening of potential peptide or small molecule inhibitors of the relevant receptor/ligand interaction. A protein or polypeptide of the present invention (including, without limitation, fragments of receptors and ligands) may be useful as inhibitors of receptor/ligand interactions. Alternatively, as described in more detail below, nucleic acids encoding proteins or polypeptides involved in receptor/ligand interactions or nucleic acids regulating the expression of such proteins or polypeptides may be introduced into appropriate host cells to increase or decrease the expression of the proteins or polypeptides as desired.

### EXAMPLE 32

#### Assaying the Proteins or Polypeptides for Anti-Inflammatory Activity

The proteins or polypeptides of the present invention may also be evaluated for anti-inflammatory activity. The anti-inflammatory activity may be achieved by providing a stimulus to cells involved in the inflammatory response, by inhibiting or promoting cell-cell interactions (such as, for example, cell adhesion), by inhibiting or promoting chemotaxis of cells involved in the inflammatory process, inhibiting or promoting cell extravasation, or by stimulating or suppressing production of other factors which more directly inhibit or promote an inflammatory response. Proteins or polypeptides exhibiting such activities can be used to treat inflammatory conditions

including chronic or acute conditions, including without limitation inflammation associated with infection (such as septic shock, sepsis or systemic inflammatory response syndrome), ischemia-reperfusion injury, endotoxin lethality, arthritis, complement-mediated hyperacute rejection, nephritis, cytokine- or chemokine-induced lung injury, inflammatory bowel disease, Crohn's disease or

5 resulting from over production of cytokines such as TNF or IL-1. Proteins or polypeptides of the invention may also be useful to treat anaphylaxis and hypersensitivity to an antigenic substance or material. Alternatively, as described in more detail below, nucleic acids encoding anti-inflammatory activity proteins or polypeptides or nucleic acids regulating the expression of such proteins or polypeptides may be introduced into appropriate host cells to increase or decrease the expression of the

10 proteins or polypeptides as desired.

### EXAMPLE 33

#### Assaying the Expressed Proteins or Polypeptides for Tumor Inhibition Activity

The proteins or polypeptides of the present invention may also be evaluated for tumor inhibition

15 activity. In addition to the activities described above for immunological treatment or prevention of tumors, a protein or polypeptide of the invention may exhibit other anti-tumor activities. A protein or polypeptide may inhibit tumor growth directly or indirectly (such as, for example, via ADCC). A protein or polypeptide may exhibit its tumor inhibitory activity by acting on tumor tissue or tumor precursor tissue, by inhibiting formation of tissues necessary to support tumor growth (such as, for

20 example, by inhibiting angiogenesis), by causing production of other factors, agents or cell types which inhibit tumor growth, or by suppressing, eliminating or inhibiting factors, agents or cell types which promote tumor growth. . Alternatively, as described in more detail below, nucleic acids encoding proteins or polypeptides with tumor inhibition activity or nucleic acids regulating the expression of such proteins or polypeptides may be introduced into appropriate host cells to increase or decrease the

25 expression of the proteins or polypeptides as desired.

A protein or polypeptide of the invention may also exhibit one or more of the following additional activities or effects: inhibiting the growth, infection or function of, or killing, infectious agents, including, without limitation, bacteria, viruses, fungi and other parasites; effecting (suppressing or enhancing) bodily characteristics, including, without limitation, height, weight, hair color, eye color,

30 skin, fat to lean ratio or other tissue pigmentation, or organ or body part size or shape (such as, for example, breast augmentation or diminution, change in bone form or shape); effecting biorhythms or circadian cycles or rhythms; effecting the fertility of male or female subjects; effecting the metabolism, catabolism, anabolism, processing, utilization, storage or elimination of dietary fat, lipid, protein, carbohydrate, vitamins, minerals, cofactors or other nutritional factors or component(s); effecting

35 behavioral characteristics, including, without limitation, appetite, libido, stress, cognition (including cognitive disorders), depression (including depressive disorders) and violent behaviors; providing analgesic effects or other pain reducing effects; promoting differentiation and growth of embryonic stem

cells in lineages other than hematopoietic lineages; hormonal or endocrine activity; in the case of enzymes, correcting deficiencies of the enzyme and treating deficiency-related diseases; treatment of hyperproliferative disorders (such as, for example, psoriasis); immunoglobulin-like activity (such as, for example, the ability to bind antigens or complement); and the ability to act as an antigen in a vaccine composition to raise an immune response against such protein or another material or entity which is cross-reactive with such protein. Alternatively, as described in more detail below, nucleic acids encoding proteins or polypeptides involved in any of the above mentioned activities or nucleic acids regulating the expression of such proteins may be introduced into appropriate host cells to increase or decrease the expression of the proteins or polypeptides as desired.

10

### EXAMPLE 34

#### Identification of Proteins or Polypeptides which Interact with Proteins or Polypeptides of the Present Invention

Proteins or polypeptides which interact with the proteins or polypeptides of the present invention, such as receptor proteins, may be identified using two hybrid systems such as the Matchmaker Two Hybrid System 2 (Catalog No. K1604-1, Clontech). As described in the manual accompanying the kit, nucleic acids encoding the proteins or polypeptides of the present invention, are inserted into an expression vector such that they are in frame with DNA encoding the DNA binding domain of the yeast transcriptional activator GAL4. cDNAs in a cDNA library which encode proteins or polypeptides which might interact with the proteins or polypeptides of the present invention are inserted into a second expression vector such that they are in frame with DNA encoding the activation domain of GAL4. The two expression plasmids are transformed into yeast and the yeast are plated on selection medium which selects for expression of selectable markers on each of the expression vectors as well as GAL4 dependent expression of the HIS3 gene. Transformants capable of growing on medium lacking histidine are screened for GAL4 dependent lacZ expression. Those cells which are positive in both the histidine selection and the lacZ assay contain plasmids encoding proteins or polypeptides which interact with the proteins or polypeptides of the present invention.

Alternatively, the system described in Lustig *et al.*, *Methods in Enzymology* 283: 83-99 (1997) may be used for identifying molecules which interact with the proteins or polypeptides of the present invention. In such systems, *in vitro* transcription reactions are performed on a pool of vectors containing nucleic acid inserts which encode the proteins or polypeptides of the present invention. The nucleic acid inserts are cloned downstream of a promoter which drives *in vitro* transcription. The resulting pools of mRNAs are introduced into *Xenopus laevis* oocytes. The oocytes are then assayed for a desired activity.

Alternatively, the pooled *in vitro* transcription products produced as described above may be translated *in vitro*. The pooled *in vitro* translation products can be assayed for a desired activity or for interaction with a known protein or polypeptide.



Proteins, polypeptides or other molecules interacting with proteins or polypeptides of the present invention can be found by a variety of additional techniques. In one method, affinity columns containing the protein or polypeptide of the present invention can be constructed. In some versions, of this method the affinity column contains chimeric proteins in which the protein or polypeptide of the present invention is fused to glutathione S-transferase. A mixture of cellular proteins or pool of expressed proteins as described above and is applied to the affinity column. Molecules interacting with the protein or polypeptide attached to the column can then be isolated and analyzed on 2-D electrophoresis gel as described in Ramunsen *et al. Electrophoresis*, **18**, 588-598 (1997). Alternatively, the molecules retained on the affinity column can be purified by electrophoresis based methods and sequenced. The same method can be used to isolate antibodies, to screen phage display products, or to screen phage display human antibodies.

Molecules interacting with the proteins or polypeptides of the present invention can also be screened by using an Optical Biosensor as described in Edwards & Leatherbarrow, *Analytical Biochemistry*, **246**, 1-6 (1997). The main advantage of the method is that it allows the determination of the association rate between the protein or polypeptide and other interacting molecules. Thus, it is possible to specifically select interacting molecules with a high or low association rate. Typically a target molecule is linked to the sensor surface (through a carboxymethyl dextran matrix) and a sample of test molecules is placed in contact with the target molecules. The binding of a test molecule to the target molecule causes a change in the refractive index and/ or thickness. This change is detected by the Biosensor provided it occurs in the evanescent field (which extends a few hundred nanometers from the sensor surface). In these screening assays, the target molecule can be one of the proteins or polypeptides of the present invention and the test sample can be a collection of proteins, polypeptides or other molecules extracted from tissues or cells, a pool of expressed proteins, combinatorial peptide and/ or chemical libraries, or phage displayed peptides. The tissues or cells from which the test molecules are extracted can originate from any species.

In other methods, a target protein or polypeptide is immobilized and the test population is a collection of unique proteins or polypeptides of the present invention.

To study the interaction of the proteins or polypeptides of the present invention with drugs, the microdialysis coupled to HPLC method described by Wang *et al.*, *Chromatographia*, **44**, 205-208(1997) or the affinity capillary electrophoresis method described by Busch *et al.*, *J. Chromatogr.* **777**:311-328 (1997) can be used.

The system described in U.S. Patent No. 5,654,150 may also be used to identify molecules which interact with the proteins or polypeptides of the present invention. In this system, pools of nucleic acids encoding the proteins or polypeptides of the present invention are transcribed and translated *in vitro* and the reaction products are assayed for interaction with a known polypeptide or antibody.

It will be appreciated by those skilled in the art that the proteins or polypeptides of the present invention may be assayed for numerous activities in addition to those specifically enumerated above.

For example, the expressed proteins or polypeptides may be evaluated for applications involving control and regulation of inflammation, tumor proliferation or metastasis, infection, or other clinical conditions. In addition, the proteins or polypeptides may be useful as nutritional agents or cosmetic agents.

The proteins or polypeptides of the present invention may be used to generate antibodies  
5 capable of specifically binding to the proteins or polypeptides of the present invention. The antibodies may be monoclonal antibodies or polyclonal antibodies. As used herein, "antibody" refers to a polypeptide or group of polypeptides which are comprised of at least one binding domain, where a binding domain is formed from the folding of variable domains of an antibody molecule to form three-dimensional binding spaces with an internal surface shape and charge distribution  
10 complementary to the features of an antigenic determinant of an antigen, which allows an immunological reaction with the antigen. Antibodies include recombinant proteins comprising the binding domains, as well as fragments, including Fab, Fab', F(ab)<sub>2</sub>, and F(ab')<sub>2</sub> fragments.

As used herein, an "antigenic determinant" is the portion of an antigen molecule, that determines the specificity of the antigen-antibody reaction. An "epitope" refers to an antigenic  
15 determinant of a polypeptide. An epitope can comprise as few as 3 amino acids in a spatial conformation which is unique to the epitope. Generally an epitope consists of at least 6 such amino acids, and more usually at least 8-10 such amino acids. Methods for determining the amino acids which make up an epitope include x-ray crystallography, 2-dimensional nuclear magnetic resonance, and epitope mapping e.g. the Pepsan method described by H. Mario Geysen *et al.* 1984. Proc. Natl.  
20 Acad. Sci. U.S.A. 81:3998-4002; PCT Publication No. WO 84/03564; and PCT Publication No. WO 84/03506.

In some embodiments, the antibodies may be capable of specifically binding to a protein or polypeptide encoded by EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids.  
25 In some embodiments, the antibody may be capable of binding an antigenic determinant or an epitope in a protein or polypeptide encoded by EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids.

In other embodiments, the antibodies may be capable of specifically binding to an EST-related  
30 polypeptide, fragment of an EST-related polypeptide, positional segment of an EST-related polypeptide or fragment of a positional segment of an EST-related polypeptide. In some embodiments, the antibody may be capable of binding an antigenic determinant or an epitope in an EST-related polypeptide, fragment of an EST-related polypeptide, positional segment of an EST-related polypeptide or fragment of a positional segment of an EST-related polypeptide.

35 In the case of secreted proteins, the antibodies may be capable of binding a full-length protein encoded by a nucleic acid of the present invention, a mature protein (*i.e.* the protein generated by

cleavage of the signal peptide) encoded by a nucleic acid of the present invention, or a signal peptide encoded by a nucleic acid of the present invention.

### EXAMPLE 35

#### 5                    Production of an Antibody to a Human Polypeptide or Protein

The above described EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids or nucleic acids encoding EST-related polypeptides, fragments of EST-related polypeptides, positional segments of EST-related polypeptides or fragments of positional segments of  
10 EST-related polypeptides are operably linked to promoters and introduced into cells as described above.

In the case of secreted proteins, nucleic acids encoding the full protein (*i.e.* the mature protein and the signal peptide), nucleic acids encoding the mature protein (*i.e.* the protein generated by cleavage of the signal peptide), or nucleic acids encoding the signal peptide are operably linked to promoters and introduced into cells as described above.

15            The encoded proteins or polypeptides are then substantially purified or isolated as described above. The concentration of protein in the final preparation is adjusted, for example, by concentration on an Amicon filter device, to the level of a few  $\mu\text{g/ml}$ . Monoclonal or polyclonal antibody to the protein or polypeptide can then be prepared as follows:

#### 1. Monoclonal Antibody Production by Hybridoma Fusion

20            Monoclonal antibody to epitopes of any of the proteins or polypeptides identified and isolated as described can be prepared from murine hybridomas according to the classical method of Kohler, and Milstein, *Nature* 256:495 (1975) or derivative methods thereof. Briefly, a mouse is repetitively inoculated with a few micrograms of the selected protein or peptides derived therefrom over a period of a few weeks. The mouse is then sacrificed, and the antibody producing cells of the spleen isolated. The  
25 spleen cells are fused by means of polyethylene glycol with mouse myeloma cells, and the excess unfused cells destroyed by growth of the system on selective media comprising aminopterin (HAT media). The successfully fused cells are diluted and aliquots of the dilution placed in wells of a microtiter plate where growth of the culture is continued. Antibody-producing clones are identified by detection of antibody in the supernatant fluid of the wells by immunoassay procedures, such as Elisa, as  
30 originally described by Engvall, *Meth. Enzymol.* 70:419 (1980). Selected positive clones can be expanded and their monoclonal antibody product harvested for use. Detailed procedures for monoclonal antibody production are described in Davis, L. *et al.* in *Basic Methods in Molecular Biology* Elsevier, New York. Section 21-2.

#### 2. Polyclonal Antibody Production by Immunization

35            Polyclonal antiserum containing antibodies to heterogenous epitopes of a single protein or polypeptide can be prepared by immunizing suitable animals with the expressed protein or peptides derived therefrom, which can be unmodified or modified to enhance immunogenicity. Effective

polyclonal antibody production is affected by many factors related both to the antigen and the host species. For example, small molecules tend to be less immunogenic than others and may require the use of carriers and adjuvant. Also, host animals response vary depending on site of inoculations and doses, with both inadequate or excessive doses of antigen resulting in low titer antisera. Small doses (ng level) of antigen administered at multiple intradermal sites appears to be most reliable. An effective immunization protocol for rabbits can be found in Vaitukaitis, *et al.* *J. Clin. Endocrinol. Metab.* 33:988-991 (1971).

Booster injections can be given at regular intervals, and antiserum harvested when antibody titer thereof, as determined semi-quantitatively, for example, by double immunodiffusion in agar against known concentrations of the antigen, begins to fall. See, for example, Ouchterlony, *et al.*, Chap. 19 in: *Handbook of Experimental Immunology* D. Wier (ed) Blackwell (1973). Plateau concentration of antibody is usually in the range of 0.1 to 0.2 mg/ml of serum (about 12  $\mu$ M). Affinity of the antisera for the antigen is determined by preparing competitive binding curves, as described, for example, by Fisher, D., Chap. 42 in: *Manual of Clinical Immunology*, 2d Ed. (Rose and Friedman, Eds.) Amer. Soc. For Microbiol., Washington, D.C. (1980).

Antibody preparations prepared according to either of the above protocols are useful in a variety of contexts. In particular, the antibodies may be used in immunoaffinity chromatography techniques such as those described below to facilitate large scale isolation, purification, or enrichment of the proteins or polypeptides encoded by EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids or for the isolation, purification or enrichment of EST-related polypeptides, fragments of EST-related polypeptides, positional segments of EST-related polypeptides or fragments of positional segments of EST-related polypeptides.

In the case of secreted proteins, the antibodies may be used for the isolation, purification, or enrichment of the full protein (*i.e.* the mature protein and the signal peptide), the mature protein (*i.e.* the protein generated by cleavage of the signal peptide), or the signal peptide are operably linked to promoters and introduced into cells as described above.

Additionally, the antibodies may be used in immunoaffinity chromatography techniques such as those described below to isolate, purify, or enrich polypeptides which have been linked to the proteins or polypeptides encoded by EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids or to isolate, purify, or enrich EST-related polypeptides, fragments of EST-related polypeptides, positional segments of EST-related polypeptides or fragments of positional segments of EST-related polypeptides.

The antibodies may also be used to determine the cellular localization of polypeptides encoded by the proteins or polypeptides encoded by EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids or the cellular

localization of EST-related polypeptides, fragments of EST-related polypeptides, positional segments of EST-related polypeptides or fragments of positional segments of EST-related polypeptides.

In addition, the antibodies may also be used to determine the cellular localization of polypeptides which have been linked to the proteins or polypeptides encoded by EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids or polypeptides which have been linked to EST-related polypeptides, fragments of EST-related polypeptides, positional segments of EST-related polypeptides or fragments of positional segments of EST-related polypeptides .

The antibodies may also be used in quantitative immunoassays which determine concentrations of antigen-bearing substances in biological samples; they may also used semi-quantitatively or qualitatively to identify the presence of antigen in a biological sample or to identify the type of tissue present in a biological sample. The antibodies may also be used in therapeutic compositions for killing cells expressing the protein or reducing the levels of the protein in the body.

#### VI. Use of 5'ESTs or Consensus Contigated 5' ESTs or Sequences Obtainable Therefrom or Portions Thereof as Reagents

The EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may be used as reagents in isolation procedures, diagnostic assays, and forensic procedures. For example, sequences from the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids, may be detectably labeled and used as probes to isolate other sequences capable of hybridizing to them. In addition, the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may be used to design PCR primers to be used in isolation, diagnostic, or forensic procedures.

##### 1. Use of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids in isolation, diagnostic and forensic procedures

#### EXAMPLE 36

##### Preparation of PCR Primers and Amplification of DNA

The EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may be used to prepare PCR primers for a variety of applications, including isolation procedures for cloning nucleic acids capable of hybridizing to such sequences, diagnostic techniques and forensic techniques. In some embodiments, the PCR primers at least 10, 15, 18, 20, 23, 25, 28, 30, 40, or 50 nucleotides in length. In some embodiments, the PCR primers may be more than 30 bases in length. It is preferred that the primer pairs have approximately the same G/C ratio, so that melting temperatures are approximately the same. A variety of PCR techniques are familiar to those skilled in the art. For a review of PCR technology, see Molecular Cloning to

Genetic Engineering White, B.A. Ed. in *Methods in Molecular Biology* 67: Humana Press, Totowa 1997.

In each of these PCR procedures, PCR primers on either side of the nucleic acid sequences to be amplified are added to a suitably prepared nucleic acid sample along with dNTPs and a thermostable polymerase such as Taq polymerase, Pfu polymerase, or Vent polymerase. The nucleic acid in the sample is denatured and the PCR primers are specifically hybridized to complementary nucleic acid sequences in the sample. The hybridized primers are extended. Thereafter, another cycle of denaturation, hybridization, and extension is initiated. The cycles are repeated multiple times to produce an amplified fragment containing the nucleic acid sequence between the primer sites.

10

### EXAMPLE 37

#### Use of the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids as probes

Probes derived from EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may be labeled with detectable labels familiar to those skilled in the art, including radioisotopes and non-radioactive labels, to provide a detectable probe. The detectable probe may be single stranded or double stranded and may be made using techniques known in the art, including *in vitro* transcription, nick translation, or kinase reactions. A nucleic acid sample containing a sequence capable of hybridizing to the labeled probe is contacted with the labeled probe. If the nucleic acid in the sample is double stranded, it may be denatured prior to contacting the probe. In some applications, the nucleic acid sample may be immobilized on a surface such as a nitrocellulose or nylon membrane. The nucleic acid sample may comprise nucleic acids obtained from a variety of sources, including genomic DNA, cDNA libraries, RNA, or tissue samples.

Procedures used to detect the presence of nucleic acids capable of hybridizing to the detectable probe include well known techniques such as Southern blotting, Northern blotting, dot blotting, colony hybridization, and plaque hybridization. In some applications, the nucleic acid capable of hybridizing to the labeled probe may be cloned into vectors such as expression vectors, sequencing vectors, or *in vitro* transcription vectors to facilitate the characterization and expression of the hybridizing nucleic acids in the sample. For example, such techniques may be used to isolate and clone sequences in a genomic library or cDNA library which are capable of hybridizing to the detectable probe as described in Example 20 above.

PCR primers made as described in Example 36 above may be used in forensic analyses, such as the DNA fingerprinting techniques described in Examples 38-42 below. Such analyses may utilize detectable probes or primers based on the sequences of the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids.

35

### EXAMPLE 38

#### Forensic Matching by DNA Sequencing

In one exemplary method, DNA samples are isolated from forensic specimens of, for example, hair, semen, blood or skin cells by conventional methods. A panel of PCR primers based on a number of the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids is then utilized in accordance with Example 36 to  
5 amplify DNA of approximately 100-200 bases in length from the forensic specimen. Corresponding sequences are obtained from a test subject. Each of these identification DNAs is then sequenced using standard techniques, and a simple database comparison determines the differences, if any, between the sequences from the subject and those from the sample. Statistically significant differences between the suspect's DNA sequences and those from the sample conclusively prove a lack of identity. This lack of  
10 identity can be proven, for example, with only one sequence. Identity, on the other hand, should be demonstrated with a large number of sequences, all matching. Preferably, a minimum of 50 statistically identical sequences of 100 bases in length are used to prove identity between the suspect and the sample.

#### EXAMPLE 39

##### 15 Positive Identification by DNA Sequencing

The technique outlined in the previous example may also be used on a larger scale to provide a unique fingerprint-type identification of any individual. In this technique, primers are prepared from a large number of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids. Preferably, 20 to 50 different primers are  
20 used. These primers are used to obtain a corresponding number of PCR-generated DNA segments from the individual in question in accordance with Example 34. Each of these DNA segments is sequenced, using the methods set forth in Example 36. The database of sequences generated through this procedure uniquely identifies the individual from whom the sequences were obtained. The same panel of primers may then be used at any later time to absolutely correlate tissue or other biological specimen with that  
25 individual.

#### EXAMPLE 40

##### Southern Blot Forensic Identification

The procedure of Example 38 is repeated to obtain a panel of at least 10 amplified sequences  
30 from an individual and a specimen. Preferably, the panel contains at least 50 amplified sequences. More preferably, the panel contains 100 amplified sequences. In some embodiments, the panel contains 200 amplified sequences. This PCR-generated DNA is then digested with one or a combination of, preferably, four base specific restriction enzymes. Such enzymes are commercially available and known to those of skill in the art. After digestion, the resultant gene fragments are size separated in multiple  
35 duplicate wells on an agarose gel and transferred to nitrocellulose using Southern blotting techniques well known to those with skill in the art. For a review of Southern blotting see Davis *et al.* (Basic Methods in Molecular Biology, 1986, Elsevier Press. pp 62-65).

A panel of probes based on the sequences of the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids are radioactively or colorimetrically labeled using methods known in the art, such as nick translation or end labeling, and hybridized to the Southern blot using techniques known in the art (Davis *et al.*, supra).

5 Preferably, the probe is at least 10, 12, 15, 18, 20, 25, 28, 30, 35, 40, 50, 75, 100, 150, 200, 300, 400 or 500 nucleotides in length. Preferably, the probes are at least 10, 12, 15, 18, 20, 25, 28, 30, 35, 40, 50, 75, 100, 150, 200, 300, 400 or 500 nucleotides in length. In some embodiments, the probes are oligonucleotides which are 40 nucleotides in length or less.

Preferably, at least 5 to 10 of these labeled probes are used, and more preferably at least about  
10 20 or 30 are used to provide a unique pattern. The resultant bands appearing from the hybridization of a large sample of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids will be a unique identifier. Since the restriction enzyme cleavage will be different for every individual, the band pattern on the Southern blot will also be unique. Increasing the number of probes will provide a statistically higher level of confidence in the  
15 identification since there will be an increased number of sets of bands used for identification.

#### EXAMPLE 41

##### Dot Blot Identification Procedure

Another technique for identifying individuals using the EST-related nucleic acids, positional  
20 segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids disclosed herein utilizes a dot blot hybridization technique.

Genomic DNA is isolated from nuclei of subject to be identified. Probes are prepared that correspond to at least 10, preferably 50 sequences from the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids.  
25 The probes are used to hybridize to the genomic DNA through conditions known to those in the art. The oligonucleotides are end labeled with  $P^{32}$  using polynucleotide kinase (Pharmacia). Dot Blots are created by spotting the genomic DNA onto nitrocellulose or the like using a vacuum dot blot manifold (BioRad, Richmond California). The nitrocellulose filter containing the genomic sequences is baked or UV linked to the filter, prehybridized and hybridized with labeled probe using techniques known in the art (Davis *et al.*, supra). The  $^{32}P$  labeled DNA fragments are sequentially hybridized with successively stringent  
30 conditions to detect minimal differences between the 30 bp sequence and the DNA. Tetramethylammonium chloride is useful for identifying clones containing small numbers of nucleotide mismatches (Wood *et al.*, *Proc. Natl. Acad. Sci. USA* 82(6):1585-1588 (1985)). A unique pattern of dots distinguishes one individual from another individual.

35 EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids can be used as probes in the following alternative



fingerprinting technique. In some embodiments, the probes are oligonucleotides which are 40 nucleotides in length or less.

Preferably, a plurality of probes having sequences from different EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids are used in the alternative fingerprinting technique. Example 42 below provides a representative alternative fingerprinting procedure in which the probes are derived from EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids.

10

#### EXAMPLE 42

##### Alternative "Fingerprint" Identification Technique

Oligonucleotides are prepared from a large number, e.g. 50, 100, or 200, EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids using commercially available oligonucleotide services such as Genset, Paris, France. Preferably, the oligonucleotides are at least 10, 15, 18, 20, 23, 25, 28, or 30 nucleotides in length. However, in some embodiments, the oligonucleotides may be more than 40, 50, 60 or 70 nucleotides in length.

Cell samples from the test subject are processed for DNA using techniques well known to those with skill in the art. The nucleic acid is digested with restriction enzymes such as EcoRI and XbaI. Following digestion, samples are applied to wells for electrophoresis. The procedure, as known in the art, may be modified to accommodate polyacrylamide electrophoresis, however in this example, samples containing 5 ug of DNA are loaded into wells and separated on 0.8% agarose gels. The gels are transferred onto nitrocellulose using standard Southern blotting techniques.

10 ng of each of the oligonucleotides are pooled and end-labeled with  $P^{32}$ . The nitrocellulose is prehybridized with blocking solution and hybridized with the labeled probes. Following hybridization and washing, the nitrocellulose filter is exposed to X-Omat AR X-ray film. The resulting hybridization pattern will be unique for each individual.

It is additionally contemplated within this example that the number of probe sequences used can be varied for additional accuracy or clarity.

30

In addition to their applications in forensics and identification, EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may be mapped to their chromosomal locations. Example 41 below describes radiation hybrid (RH) mapping of human chromosomal regions using EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids. Example 42 below describes a representative procedure for mapping EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids to their locations on human chromosomes. Example 43 below describes mapping of

EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids on metaphase chromosomes by Fluorescence In Situ Hybridization (FISH).

- 5 2. Use of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids in Chromosome Mapping

#### EXAMPLE 43

##### Radiation hybrid mapping of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids to the human genome

10 Radiation hybrid (RH) mapping is a somatic cell genetic approach that can be used for high resolution mapping of the human genome. In this approach, cell lines containing one or more human chromosomes are lethally irradiated, breaking each chromosome into fragments whose size depends on the radiation dose. These fragments are rescued by fusion with cultured rodent cells, yielding subclones  
15 containing different portions of the human genome. This technique is described by Benham *et al.* (*Genomics* 4:509-517, 1989) and Cox *et al.*, (*Science* 250:245-250, 1990). The random and independent nature of the subclones permits efficient mapping of any human genome marker. Human DNA isolated from a panel of 80-100 cell lines provides a mapping reagent for ordering EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related  
20 nucleic acids. In this approach, the frequency of breakage between markers is used to measure distance, allowing construction of fine resolution maps as has been done using conventional ESTs (Schuler *et al.*, *Science* 274:540-546, 1996).

RH mapping has been used to generate a high-resolution whole genome radiation hybrid map of human chromosome 17q22-q25.3 across the genes for growth hormone (GH) and thymidine kinase (TK)  
25 (Foster *et al.*, *Genomics* 33:185-192, 1996), the region surrounding the Gorlin syndrome gene (Obermayr *et al.*, *Eur. J. Hum. Genet.* 4:242-245, 1996), 60 loci covering the entire short arm of chromosome 12 (Raeymaekers *et al.*, *Genomics* 29:170-178, 1995), the region of human chromosome 22 containing the neurofibromatosis type 2 locus (Frazer *et al.*, *Genomics* 14:574-584, 1992) and 13 loci on the long arm of chromosome 5 (Warrington *et al.*, *Genomics* 11:701-708, 1991).

30

#### EXAMPLE 44

##### Mapping of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids to Human Chromosomes using PCR techniques

35 EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may be assigned to human chromosomes using PCR based methodologies. In such approaches, oligonucleotide primer pairs are designed from EST-related

nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids to minimize the chance of amplifying through an intron. Preferably, the oligonucleotide primers are 18-23 bp in length and are designed for PCR amplification. The creation of PCR primers from known sequences is well known to those with skill in the art. For a review of PCR technology see Erlich, in PCR Technology; Principles and Applications for DNA Amplification. 1992. W.H. Freeman and Co., New York.

The primers are used in polymerase chain reactions (PCR) to amplify templates from total human genomic DNA. PCR conditions are as follows: 60 ng of genomic DNA is used as a template for PCR with 80 ng of each oligonucleotide primer, 0.6 unit of Taq polymerase, and 1  $\mu$ Cu of a  $^{32}$ P-labeled deoxycytidine triphosphate. The PCR is performed in a microplate thermocycler (Techne) under the following conditions: 30 cycles of 94°C, 1.4 min; 55°C, 2 min; and 72°C, 2 min; with a final extension at 72°C for 10 min. The amplified products are analyzed on a 6% polyacrylamide sequencing gel and visualized by autoradiography. If the length of the resulting PCR product is identical to the distance between the ends of the primer sequences in the 5'EST from which the primers are derived, then the PCR reaction is repeated with DNA templates from two panels of human-rodent somatic cell hybrids, BIOS PCRable DNA (BIOS Corporation) and NIGMS Human-Rodent Somatic Cell Hybrid Mapping Panel Number 1 (NIGMS, Camden, NJ).

PCR is used to screen a series of somatic cell hybrid cell lines containing defined sets of human chromosomes for the presence of a given 5'EST. DNA is isolated from the somatic hybrids and used as starting templates for PCR reactions using the primer pairs from the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids. Only those somatic cell hybrids with chromosomes containing the human gene corresponding to the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids will yield an amplified fragment. The 5'ESTs are assigned to a chromosome by analysis of the segregation pattern of PCR products from the somatic hybrid DNA templates. The single human chromosome present in all cell hybrids that give rise to an amplified fragment is the chromosome containing that EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids. For a review of techniques and analysis of results from somatic cell gene mapping experiments. (See Ledbetter *et al.*, Genomics 6:475-481 (1990)).

Alternatively, the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may be mapped to individual chromosomes using FISH as described in Example 45 below.

#### EXAMPLE 45

Mapping of EST-related nucleic acids, positional segments of  
EST-related nucleic acids or fragments of positional segments of

EST-related nucleic acids to Chromosomes Using  
Fluorescence *In Situ* Hybridization

Fluorescence in situ hybridization allows the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids to be mapped to a particular location on a given chromosome. The chromosomes to be used for fluorescence in situ hybridization techniques may be obtained from a variety of sources including cell cultures, tissues, or whole blood.

In a preferred embodiment, chromosomal localization of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids are obtained by FISH as described by Cherif *et al.* (*Proc. Natl. Acad. Sci. U.S.A.*, 87:6639-6643, 1990). Metaphase chromosomes are prepared from phytohemagglutinin (PHA)-stimulated blood cell donors. PHA-stimulated lymphocytes from healthy males are cultured for 72 h in RPMI-1640 medium. For synchronization, methotrexate (10  $\mu$ M) is added for 17 h, followed by addition of 5-bromodeoxyuridine (5-BrdU, 0.1 mM) for 6 h. Colcemid (1  $\mu$ g/ml) is added for the last 15 min before harvesting the cells. Cells are collected, washed in RPMI, incubated with a hypotonic solution of KCl (75 mM) at 37°C for 15 min and fixed in three changes of methanol:acetic acid (3:1). The cell suspension is dropped onto a glass slide and air dried. The EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids is labeled with biotin-16 dUTP by nick translation according to the manufacturer's instructions (Bethesda Research Laboratories, Bethesda, MD), purified using a Sephadex G-50 column (Pharmacia, Upsala, Sweden) and precipitated. Just prior to hybridization, the DNA pellet is dissolved in hybridization buffer (50% formamide, 2 X SSC, 10% dextran sulfate, 1 mg/ml sonicated salmon sperm DNA, pH 7) and the probe is denatured at 70°C for 5-10 min.

Slides kept at -20°C are treated for 1 h at 37°C with RNase A (100  $\mu$ g/ml), rinsed three times in 2 X SSC and dehydrated in an ethanol series. Chromosome preparations are denatured in 70% formamide, 2 X SSC for 2 min at 70°C, then dehydrated at 4°C. The slides are treated with proteinase K (10  $\mu$ g/100 ml in 20 mM Tris-HCl, 2 mM CaCl<sub>2</sub>) at 37°C for 8 min and dehydrated. The hybridization mixture containing the probe is placed on the slide, covered with a coverslip, sealed with rubber cement and incubated overnight in a humid chamber at 37°C. After hybridization and post-hybridization washes, the biotinylated probe is detected by avidin-FITC and amplified with additional layers of biotinylated goat anti-avidin and avidin-FITC. For chromosomal localization, fluorescent R-bands are obtained as previously described (Cherif *et al.*, *supra.*). The slides are observed under a LEICA fluorescence microscope (DMRXA). Chromosomes are counterstained with propidium iodide and the fluorescent signal of the probe appears as two symmetrical yellow-green spots on both chromatids of the fluorescent R-band chromosome (red). Thus, a particular EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may be localized to a particular cytogenetic R-band on a given chromosome.

Once the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids have been assigned to particular chromosomes using the techniques described in Examples 42-44 above, they may be utilized to construct a high resolution map of the chromosomes on which they are located or to identify the chromosomes in a sample.

#### EXAMPLE 46

##### Use of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids to Construct or Expand Chromosome Maps

Chromosome mapping involves assigning a given unique sequence to a particular chromosome as described above. Once the unique sequence has been mapped to a given chromosome, it is ordered relative to other unique sequences located on the same chromosome. One approach to chromosome mapping utilizes a series of yeast artificial chromosomes (YACs) bearing several thousand long inserts derived from the chromosomes of the organism from which the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids are obtained. This approach is described in Ramaiah Nagaraja *et al.*, *Genome Research* 7:210-222, March 1997. Briefly, in this approach each chromosome is broken into overlapping pieces which are inserted into the YAC vector. The YAC inserts are screened using PCR or other methods to determine whether they include the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids whose position is to be determined. Once an insert has been found which includes the 5'EST, the insert can be analyzed by PCR or other methods to determine whether the insert also contains other sequences known to be on the chromosome or in the region from which the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids was derived. This process can be repeated for each insert in the YAC library to determine the location of each of the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids relative to one another and to other known chromosomal markers. In this way, a high resolution map of the distribution of numerous unique markers along each of the organisms chromosomes may be obtained.

As described in Example 47 below EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may also be used to identify genes associated with a particular phenotype, such as hereditary disease or drug response.

##### 3. Use of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids Gene Identification

#### EXAMPLE 47

##### Identification of genes associated with hereditary diseases or drug response

This example illustrates an approach useful for the association of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids with particular phenotypic characteristics. In this example, a particular EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids is used as a test probe to associate that EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids with a particular phenotypic characteristic.

EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids are mapped to a particular location on a human chromosome using techniques such as those described in Examples 41 and 42 or other techniques known in the art. A search of Mendelian Inheritance in Man (V. McKusick, *Mendelian Inheritance in Man* (available on line through Johns Hopkins University Welch Medical Library) reveals the region of the human chromosome which contains the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids to be a very gene rich region containing several known genes and several diseases or phenotypes for which genes have not been identified. The gene corresponding to this EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids thus becomes an immediate candidate for each of these genetic diseases.

Cells from patients with these diseases or phenotypes are isolated and expanded in culture. PCR primers from the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids are used to screen genomic DNA, mRNA or cDNA obtained from the patients. EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids that are not amplified in the patients can be positively associated with a particular disease by further analysis. Alternatively, the PCR analysis may yield fragments of different lengths when the samples are derived from an individual having the phenotype associated with the disease than when the sample is derived from a healthy individual, indicating that the gene containing the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may be responsible for the genetic disease.

30

## **VII. Use of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids to Construct Vectors**

The present EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may also be used to construct secretion vectors capable of directing the secretion of the proteins encoded by genes therein. Such secretion vectors may facilitate the purification or enrichment of the proteins encoded by genes inserted therein by

35

reducing the number of background proteins from which the desired protein must be purified or enriched. Exemplary secretion vectors are described in Example 48 below.

#### 1. Construction of secretion vectors

### **EXAMPLE 48**

#### **5                   Construction of Secretion Vectors**

The secretion vectors of the present invention include a promoter capable of directing gene expression in the host cell, tissue, or organism of interest. Such promoters include the Rous Sarcoma Virus promoter, the SV40 promoter, the human cytomegalovirus promoter, and other promoters familiar to those skilled in the art.

10           A signal sequence from one of the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids is operably linked to the promoter such that the mRNA transcribed from the promoter will direct the translation of the signal peptide. Preferably, the signal sequence is from one of the nucleic acids of SEQ ID NOs.24-811. The host cell, tissue, or organism may be any cell, tissue, or organism which recognizes the signal  
15 peptide encoded by the signal sequence in the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids. Suitable hosts include mammalian cells, tissues or organisms, avian cells, tissues, or organisms, insect cells, tissues or organisms, or yeast.

In addition, the secretion vector contains cloning sites for inserting genes encoding the proteins  
20 which are to be secreted. The cloning sites facilitate the cloning of the insert gene in frame with the signal sequence such that a fusion protein in which the signal peptide is fused to the protein encoded by the inserted gene is expressed from the mRNA transcribed from the promoter. The signal peptide directs the extracellular secretion of the fusion protein.

The secretion vector may be DNA or RNA and may integrate into the chromosome of the host,  
25 be stably maintained as an extrachromosomal replicon in the host, be an artificial chromosome, or be transiently present in the host. Preferably, the secretion vector is maintained in multiple copies in each host cell. As used herein, multiple copies means at least 2, 5, 10, 20, 25, 50 or more than 50 copies per cell. In some embodiments, the multiple copies are maintained extrachromosomally. In other embodiments, the multiple copies result from amplification of a chromosomal sequence.

30           Many nucleic acid backbones suitable for use as secretion vectors are known to those skilled in the art, including retroviral vectors, SV40 vectors, Bovine Papilloma Virus vectors, yeast integrating plasmids, yeast episomal plasmids, yeast artificial chromosomes, human artificial chromosomes, P element vectors, baculovirus vectors, or bacterial plasmids capable of being transiently introduced into the host.

35           The secretion vector may also contain a polyA signal such that the polyA signal is located downstream of the gene inserted into the secretion vector.

After the gene encoding the protein for which secretion is desired is inserted into the secretion vector, the secretion vector is introduced into the host cell, tissue, or organism using calcium phosphate precipitation, DEAE-Dextran, electroporation, liposome-mediated transfection, viral particles or as naked DNA. The protein encoded by the inserted gene is then purified or enriched from the supernatant using conventional techniques such as ammonium sulfate precipitation, immunoprecipitation, immunoaffinitychromatography, size exclusion chromatography, ion exchange chromatography, and HPLC. Alternatively, the secreted protein may be in a sufficiently enriched or pure state in the supernatant or growth media of the host to permit it to be used for its intended purpose without further enrichment.

The signal sequences may also be inserted into vectors designed for gene therapy. In such vectors, the signal sequence is operably linked to a promoter such that mRNA transcribed from the promoter encodes the signal peptide. A cloning site is located downstream of the signal sequence such that a gene encoding a protein whose secretion is desired may readily be inserted into the vector and fused to the signal sequence. The vector is introduced into an appropriate host cell. The protein expressed from the promoter is secreted extracellularly, thereby producing a therapeutic effect.

#### EXAMPLE 49

##### Fusion Vectors

The EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may be used to construct fusion vectors for the expression of chimeric polypeptides. The chimeric polypeptides comprise a first polypeptide portion and a second polypeptide portion. In the fusion vectors of the present invention, nucleic acids encoding the first polypeptide portion and the second polypeptide portion are joined in frame with one another so as to generate a nucleic acid encoding the chimeric polypeptide. The nucleic acid encoding the chimeric polypeptide is operably linked to a promoter which directs the expression of an mRNA encoding the chimeric polypeptide. The promoter may be in any of the expression vectors described herein including those described in Examples 21 and 48.

Preferably, the fusion vector is maintained in multiple copies in each host cell. In some embodiments, the multiple copies are maintained extrachromosomally. In other embodiments, the multiple copies result from amplification of a chromosomal sequence.

The first polypeptide portion may comprise any of the polypeptides encoded by the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids. In some embodiments, the first polypeptide portion may be one of the EST-related polypeptides, fragments of EST-related polypeptides, positional segments of EST-related polypeptides, or fragments of positional segments of EST-related polypeptides.

The second polypeptide portion may comprise any polypeptide of interest. In some embodiments, the second polypeptide portion may comprise a polypeptide having a detectable



- enzymatic activity such as green fluorescent protein or  $\beta$  galactosidase. Chimeric polypeptides in which the second polypeptide portion comprises a detectable polypeptide may be used to determine the intracellular localization of the first polypeptide portion. In such procedures, the fusion vector encoding the chimeric polypeptide is introduced into a host cell under conditions which facilitate the expression of
- 5 the chimeric polypeptide. Where appropriate, the cells are treated with a detection reagent which is visible under the microscope following a catalytic reaction with the detectable polypeptide and the cellular location of the detection reagent is determined. For example, if the polypeptide having a detectable enzymatic activity is  $\beta$  galactosidase, the cells may be treated with Xgal. Alternatively, where the detectable polypeptide is directly detectable without the addition of a detection reagent, the
- 10 intracellular location of the chimeric polypeptide is determined by performing microscopy under conditions in which the detectable polypeptide is visible. For example, if the detectable polypeptide is green fluorescent protein or a modified version thereof, microscopy is performed by exposing the host cells to light having an appropriate wavelength to cause the green fluorescent protein or modified version thereof to fluoresce.
- 15 Alternatively, the second polypeptide portion may comprise a polypeptide whose isolation, purification, or enrichment is desired. In such embodiments, the isolation, purification, or enrichment of the second polypeptide portion may be achieved by performing the immunoaffinity chromatography procedures described below using an immunoaffinity column having an antibody directed against the first polypeptide portion coupled thereto.
- 20 The proteins encoded by the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids or the EST-related polypeptides, fragments of EST-related polypeptides, positional segments of EST-related polypeptides, or fragments of positional segments of EST-related polypeptides may also be used to generate antibodies as explained herein in order to identify the tissue type or cell species from which a sample is derived as
- 25 described in Example 50.

### EXAMPLE 50

#### Identification of Tissue Types or Cell Species by Means of Labeled Tissue Specific Antibodies

- 30 Identification of specific tissues is accomplished by the visualization of tissue specific antigens by means of antibody preparations as described herein which are conjugated, directly or indirectly to a detectable marker. Selected labeled antibody species bind to their specific antigen binding partner in tissue sections, cell suspensions, or in extracts of soluble proteins from a tissue sample to provide a pattern for qualitative or semi-qualitative interpretation.
- 35 Antisera for these procedures must have a potency exceeding that of the native preparation, and for that reason, antibodies are concentrated to a mg/ml level by isolation of the gamma globulin fraction, for example, by ion-exchange chromatography or by ammonium sulfate fractionation. Also, to provide

the most specific antisera, unwanted antibodies, for example to common proteins, must be removed from the gamma globulin fraction, for example by means of insoluble immunoabsorbents, before the antibodies are labeled with the marker. Either monoclonal or heterologous antisera is suitable for either procedure.

#### 5 1. Immunohistochemical Techniques

Purified, high-titer antibodies, prepared as described above, are conjugated to a detectable marker, as described, for example, by Fudenberg, H., Chap. 26 in: *Basic 503 Clinical Immunology*, 3<sup>rd</sup> Ed. Lange, Los Altos, California (1980) or Rose, . *et al.*, Chap. 12 in: *Methods in Immunodiagnosis*, 2d Ed. John Wiley and Sons, New York (1980).

10 A fluorescent marker, either fluorescein or rhodamine, is preferred, but antibodies can also be labeled with an enzyme that supports a color producing reaction with a substrate, such as horseradish peroxidase. Markers can be added to tissue-bound antibody in a second step, as described below. Alternatively, the specific antitissue antibodies can be labeled with ferritin or other electron dense particles, and localization of the ferritin coupled antigen-antibody complexes achieved by means of an  
15 electron microscope. In yet another approach, the antibodies are radiolabeled, with, for example <sup>125</sup>I, and detected by overlaying the antibody treated preparation with photographic emulsion.

Preparations to carry out the procedures can comprise monoclonal or polyclonal antibodies to a single protein or peptide identified as specific to a tissue type, for example, brain tissue, or antibody preparations to several antigenically distinct tissue specific antigens can be used in panels, independently  
20 or in mixtures, as required.

Tissue sections and cell suspensions are prepared for immunohistochemical examination according to common histological techniques. Multiple cryostat sections (about 4 µm, unfixed) of the unknown tissue and known control, are mounted and each slide covered with different dilutions of the antibody preparation. Sections of known and unknown tissues should also be treated with preparations  
25 to provide a positive control, a negative control, for example, pre-immune sera, and a control for non-specific staining, for example, buffer.

Treated sections are incubated in a humid chamber for 30 min at room temperature, rinsed, then washed in buffer for 30-45 min. Excess fluid is blotted away, and the marker developed.

If the tissue specific antibody was not labeled in the first incubation, it can be labeled at this time  
30 in a second antibody-antibody reaction, for example, by adding fluorescein- or enzyme-conjugated antibody against the immunoglobulin class of the antiserum-producing species, for example, fluorescein labeled antibody to mouse IgG. Such labeled sera are commercially available.

The antigen found in the tissues by the above procedure can be quantified by measuring the intensity of color or fluorescence on the tissue section, and calibrating that signal using appropriate  
35 standards.

#### 2. Identification of Tissue Specific Soluble Proteins

The visualization of tissue specific proteins and identification of unknown tissues from that procedure is carried out using the labeled antibody reagents and detection strategy as described for immunohistochemistry; however the sample is prepared according to an electrophoretic technique to distribute the proteins extracted from the tissue in an orderly array on the basis of molecular weight for  
5 detection.

A tissue sample is homogenized using a Virtis apparatus; cell suspensions are disrupted by Dounce homogenization or osmotic lysis, using detergents in either case as required to disrupt cell membranes, as is the practice in the art. Insoluble cell components such as nuclei, microsomes, and membrane fragments are removed by ultracentrifugation, and the soluble protein-containing fraction  
10 concentrated if necessary and reserved for analysis.

A sample of the soluble protein solution is resolved into individual protein species by conventional SDS polyacrylamide electrophoresis as described, for example, by Davis, L. *et al.*, Section 19-2 in: *Basic Methods in Molecular Biology* (P. Leder, ed), Elsevier, New York (1986), using a range of amounts of polyacrylamide in a set of gels to resolve the entire molecular weight range of proteins to  
15 be detected in the sample. A size marker is run in parallel for purposes of estimating molecular weights of the constituent proteins. Sample size for analysis is a convenient volume of from 5 to 55  $\mu$ l, and containing from about 1 to 100  $\mu$ g protein. An aliquot of each of the resolved proteins is transferred by blotting to a nitrocellulose filter paper, a process that maintains the pattern of resolution. Multiple copies are prepared. The procedure, known as Western Blot Analysis, is well described in Davis, L. *et al.*,  
20 *supra* Section 19-3. One set of nitrocellulose blots is stained with Coomassie Blue dye to visualize the entire set of proteins for comparison with the antibody bound proteins. The remaining nitrocellulose filters are then incubated with a solution of one or more specific antisera to tissue specific proteins prepared as described in Examples 20 and 33. In this procedure, as in procedure A above, appropriate positive and negative sample and reagent controls are run.

25 In either procedure described above a detectable label can be attached to the primary tissue antigen-primary antibody complex according to various strategies and permutations thereof. In a straightforward approach, the primary specific antibody can be labeled; alternatively, the unlabeled complex can be bound by a labeled secondary anti-IgG antibody. In other approaches, either the primary or secondary antibody is conjugated to a biotin molecule, which can, in a subsequent step, bind an avidin  
30 conjugated marker. According to yet another strategy, enzyme labeled or radioactive protein A, which has the property of binding to any IgG, is bound in a final step to either the primary or secondary antibody.

## EXAMPLE 51

### 35 Immunohistochemical Localization of Polypeptides

The antibodies prepared as described herein above may be utilized to determine the cellular location of a polypeptide. The polypeptide may be any of the polypeptides encoded by EST-related

nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids or the polypeptide may be one of the EST-related polypeptides, fragments of EST-related polypeptides, positional segments of EST-related polypeptides, or fragments of positional segments of EST-related polypeptides. In some embodiments, the polypeptide may be a chimeric  
5 polypeptide such as those encoded by the fusion vectors of Example 49.

Cells expressing the polypeptide to be localized are applied to a microscope slide and fixed using any of the procedures typically employed in immunohistochemical localization techniques, including the methods described in *Current Protocols in Molecular Biology*, John Wiley and Sons, Inc. 1997. Following a washing step, the cells are contacted with the antibody. In some embodiments, the  
10 antibody is conjugated to a detectable marker as described above to facilitate detection. Alternatively, in some embodiments, after the cells have been contacted with an antibody to the polypeptide to be localized, a secondary antibody which has been conjugated to a detectable marker is placed in contact with the antibody against the polypeptide to be localized.

Thereafter, microscopy is performed under conditions suitable for visualizing the cellular  
15 location of the polypeptide.

The visualization of tissue specific antigen binding at levels above those seen in control tissues to one or more tissue specific antibodies, directed against the polypeptides encoded by EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids or antibodies against the EST-related polypeptides, fragments of EST-related  
20 polypeptides, positional segments of EST-related polypeptides, or fragments of positional segments of EST-related polypeptides, can identify tissues of unknown origin, for example, forensic samples, or differentiated tumor tissue that has metastasized to foreign bodily sites.

The antibodies described herein may also be used in the immunoaffinity chromatography techniques described below to isolate, purify or enrich the polypeptides encoded by the EST-related  
25 nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids or to isolate, purify or enrich EST-related polypeptides, fragments of EST-related polypeptides, positional segments of EST-related polypeptides, or fragments of positional segments of EST-related polypeptides. The immunoaffinity chromatography techniques described below may also be used to isolate, purify or enrich polypeptides which have been linked to the  
30 polypeptides encoded by the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids or to isolate, purify or enrich polypeptides which have been linked to EST-related polypeptides, fragments of EST-related polypeptides, positional segments of EST-related polypeptides, or fragments of positional segments of EST-related polypeptides.

**EXAMPLE 52**Immunoaffinity Chromatography

Antibodies prepared as described above are coupled to a support. Preferably, the antibodies are monoclonal antibodies, but polyclonal antibodies may also be used. The support may be any of those typically employed in immunoaffinity chromatography, including Sepharose CL-4B (Pharmacia, Piscataway, NJ), Sepharose CL-2B (Pharmacia, Piscataway, NJ), Affi-gel 10 (Biorad, Richmond, CA), or glass beads.

The antibodies may be coupled to the support using any of the coupling reagents typically used in immunoaffinity chromatography, including cyanogen bromide. After coupling the antibody to the support, the support is contacted with a sample which contains a target polypeptide whose isolation, purification or enrichment is desired. The target polypeptide may be a polypeptide encoded by the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids or the target polypeptide may be one of the EST-related polypeptides, fragments of EST-related polypeptides, positional segments of EST-related polypeptides, or fragments of positional segments of EST-related polypeptides. The target polypeptides may also be polypeptides which have been linked to the polypeptides encoded by the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids or the target polypeptides may be polypeptides which have been linked to EST-related polypeptides, fragments of EST-related polypeptides, positional segments of EST-related polypeptides, or fragments of positional segments of EST-related polypeptides using the fusion vectors described above.

Preferably, the sample is placed in contact with the support for a sufficient amount of time and under appropriate conditions to allow at least 50% of the target polypeptide to specifically bind to the antibody coupled to the support.

Thereafter, the support is washed with an appropriate wash solution to remove polypeptides which have non-specifically adhered to the support. The wash solution may be any of those typically employed in immunoaffinity chromatography, including PBS, Tris-lithium chloride buffer (0.1M lysine base and 0.5M lithium chloride, pH 8.0), Tris-hydrochloride buffer (0.05M Tris-hydrochloride, pH 8.0), or Tris/Triton/NaCl buffer (50mM Tris.cl, pH 8.0 or 9.0, 0.1% Triton X-100, and 0.5MNaCl).

After washing, the specifically bound target polypeptide is eluted from the support using the high pH or low pH elution solutions typically employed in immunoaffinity chromatography. In particular, the elution solutions may contain an eluant such as triethanolamine, diethylamine, calcium chloride, sodium thiocyanate, potassium bromide, acetic acid, or glycine. In some embodiments, the elution solution may also contain a detergent such as Triton X-100 or octyl- $\beta$ -D-glucoside.

The EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may also be used to clone sequences located upstream of the 5'ESTs which are capable of regulating gene expression, including promoter sequences, enhancer

sequences, and other upstream sequences which influence transcription or translation levels. Once identified and cloned, these upstream regulatory sequences may be used in expression vectors designed to direct the expression of an inserted gene in a desired spatial, temporal, developmental, or quantitative fashion. Example 51 describes a method for cloning sequences upstream of the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids.

## 2. Identification of upstream sequences with promoting or regulatory activities

### **EXAMPLE 53**

#### 10 Use of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids to Clone Upstream Sequences from Genomic DNA

Sequences derived from EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may be used to isolate the promoters of the corresponding genes using chromosome walking techniques. In one chromosome walking technique, which utilizes the GenomeWalker™ kit available from Clontech, five complete genomic DNA samples are each digested with a different restriction enzyme which has a 6 base recognition site and leaves a blunt end. Following digestion, oligonucleotide adapters are ligated to each end of the resulting genomic DNA fragments.

For each of the five genomic DNA libraries, a first PCR reaction is performed according to the manufacturer's instructions using an outer adapter primer provided in the kit and an outer gene specific primer. The gene specific primer should be selected to be specific for 5' EST of interest and should have a melting temperature, length, and location in the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids which is consistent with its use in PCR reactions. Each first PCR reaction contains 5ng of genomic DNA, 5 µl of 10X Tth reaction buffer, 0.2 mM of each dNTP, 0.2 µM each of outer adapter primer and outer gene specific primer, 1.1 mM of Mg(OAc)<sub>2</sub>, and 1 µl of the Tth polymerase 50X mix in a total volume of 50 µl. The reaction cycle for the first PCR reaction is as follows: 1 min at 94°C / 2 sec at 94°C, 3 min at 72°C (7 cycles) / 2 sec at 94°C, 3 min at 67°C (32 cycles) / 5 min at 67°C.

The product of the first PCR reaction is diluted and used as a template for a second PCR reaction according to the manufacturer's instructions using a pair of nested primers which are located internally on the amplicon resulting from the first PCR reaction. For example, 5 µl of the reaction product of the first PCR reaction mixture may be diluted 180 times. Reactions are made in a 50 µl volume having a composition identical to that of the first PCR reaction except the nested primers are used. The first nested primer is specific for the adapter, and is provided with the GenomeWalker™ kit. The second nested primer is specific for the particular EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids for which the promoter is to be cloned and should have a melting temperature, length, and location in

the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids which is consistent with its use in PCR reactions. The reaction parameters of the second PCR reaction are as follows: 1 min at 94°C / 2 sec at 94°C, 3 min at 72°C (6 cycles) / 2 sec at 94°C, 3 min at 67°C (25 cycles) / 5 min at 67°C. The product of the second PCR reaction is purified, cloned, and sequenced using standard techniques.

Alternatively, two or more human genomic DNA libraries can be constructed by using two or more restriction enzymes. The digested genomic DNA is cloned into vectors which can be converted into single stranded, circular, or linear DNA. A biotinylated oligonucleotide comprising at least 10, 12, 15, 18, 20, 23, 25, 27, 30, 35, 40, or 50 nucleotides from the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids sequence is hybridized to the single stranded DNA. Hybrids between the biotinylated oligonucleotide and the single stranded DNA containing the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids are isolated as described above. Thereafter, the single stranded DNA containing the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids is released from the beads and converted into double stranded DNA using a primer specific for the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids or a primer corresponding to a sequence included in the cloning vector. The resulting double stranded DNA is transformed into bacteria. cDNAs containing the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids are identified by colony PCR or colony hybridization.

Once the upstream genomic sequences have been cloned and sequenced as described above, prospective promoters and transcription start sites within the upstream sequences may be identified by comparing the sequences upstream of the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids with databases containing known transcription start sites, transcription factor binding sites, or promoter sequences.

In addition, promoters in the upstream sequences may be identified using promoter reporter vectors as described in Example 54.

#### EXAMPLE 54

##### Identification of Promoters in Cloned Upstream Sequences

The genomic sequences upstream of the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids are cloned into a suitable promoter reporter vector, such as the pSEAP-Basic, pSEAP-Enhancer, p $\beta$ -gal-Basic, p $\beta$ -gal-Enhancer, or pEGFP-1 Promoter Reporter vectors available from Clontech. Briefly, each of these promoter reporter vectors include multiple cloning sites positioned upstream of a reporter gene encoding a readily assayable protein such as secreted alkaline phosphatase,  $\beta$ -galactosidase, or green fluorescent

protein. The sequences upstream of the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids are inserted into the cloning sites upstream of the reporter gene in both orientations and introduced into an appropriate host cell. The level of reporter protein is assayed and compared to the level obtained from a vector which  
5 lacks an insert in the cloning site. The presence of an elevated expression level in the vector containing the insert with respect to the control vector indicates the presence of a promoter in the insert. If necessary, the upstream sequences can be cloned into vectors which contain an enhancer for augmenting transcription levels from weak promoter sequences. A significant level of expression above that observed with the vector lacking an insert indicates that a promoter sequence is present in the inserted  
10 upstream sequence.

Appropriate host cells for the promoter reporter vectors may be chosen based on the results of the above described determination of expression patterns of the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids. For example, if the expression pattern analysis indicates that the mRNA corresponding to a particular  
15 EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids is expressed in fibroblasts, the promoter reporter vector may be introduced into a human fibroblast cell line.

Promoter sequences within the upstream genomic DNA may be further defined by constructing nested deletions in the upstream DNA using conventional techniques such as Exonuclease III digestion.  
20 The resulting deletion fragments can be inserted into the promoter reporter vector to determine whether the deletion has reduced or obliterated promoter activity. In this way, the boundaries of the promoters may be defined. If desired, potential individual regulatory sites within the promoter may be identified using site directed mutagenesis or linker scanning to obliterate potential transcription factor binding sites within the promoter individually or in combination. The effects of these mutations on transcription  
25 levels may be determined by inserting the mutations into the cloning sites in the promoter reporter vectors.

### EXAMPLE 55

#### Cloning and Identification of Promoters

30 Using the method described in Example 54 above with 5' ESTs, sequences upstream of several genes were obtained. Using the primer pairs GGG AAG ATG GAG ATA GTA TTG CCT G (SEQ ID NO:15) and CTG CCA TGT ACA TGA TAG AGA GAT TC (SEQ ID NO:16), the promoter having the internal designation P13H2 (SEQ ID NO:17) was obtained.

Using the primer pairs GTA CCA GGGG ACT GTG ACC ATT GC (SEQ ID NO:18) and CTG  
35 TGA CCA TTG CTC CCA AGA GAG (SEQ ID NO:19), the promoter having the internal designation P15B4 (SEQ ID NO:20) was obtained.



Using the primer pairs CTG GGA TGG AAG GCA CGG TA (SEQ ID NO:21) and GAG ACC ACA CAG CTA GAC AA (SEQ ID NO:22), the promoter having the internal designation P29B6 (SEQ ID NO:23) was obtained.

Figure 4 provides a schematic description of the promoters isolated and the way they are assembled with the corresponding 5' tags. The upstream sequences were screened for the presence of motifs resembling transcription factor binding sites or known transcription start sites using the computer program MatInspector release 2.0, August 1996.

Figure 5 describes the transcription factor binding sites present in each of these promoters. The columns labeled matrice provides the name of the MatInspector matrix used. The column labeled position provides the 5' position of the promoter site. Numeration of the sequence starts from the transcription site as determined by matching the genomic sequence with the 5' EST sequence. The column labeled "orientation" indicates the DNA strand on which the site is found, with the + strand being the coding strand as determined by matching the genomic sequence with the sequence of the 5' EST. The column labeled "score" provides the MatInspector score found for this site. The column labeled "length" provides the length of the site in nucleotides. The column labeled "sequence" provides the sequence of the site found.

Bacterial clones containing plasmids containing the promoter sequences described above described above are presently stored in the inventor's laboratories under the internal identification numbers provided above. The inserts may be recovered from the deposited materials by growing an aliquot of the appropriate bacterial clone in the appropriate medium. The plasmid DNA can then be isolated using plasmid isolation procedures familiar to those skilled in the art such as alkaline lysis minipreps or large scale alkaline lysis plasmid isolation procedures. If desired the plasmid DNA may be further enriched by centrifugation on a cesium chloride gradient, size exclusion chromatography, or anion exchange chromatography. The plasmid DNA obtained using these procedures may then be manipulated using standard cloning techniques familiar to those skilled in the art. Alternatively, a PCR can be done with primers designed at both ends of the inserted EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids. The PCR product which corresponds to the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids can then be manipulated using standard cloning techniques familiar to those skilled in the art.

The promoters and other regulatory sequences located upstream of the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may be used to design expression vectors capable of directing the expression of an inserted gene in a desired spatial, temporal, developmental, or quantitative manner. A promoter capable of directing the desired spatial, temporal, developmental, and quantitative patterns may be selected using the results of the expression analysis described above. For example, if a promoter which confers a high level of expression in muscle is desired, the promoter sequence upstream of EST-related nucleic acids,

positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids derived from an mRNA which are expressed at a high level in muscle, as determined by the methods above, may be used in the expression vector.

Preferably, the desired promoter is placed near multiple restriction sites to facilitate the cloning of the desired insert downstream of the promoter, such that the promoter is able to drive expression of the inserted gene. The promoter may be inserted in conventional nucleic acid backbones designed for extrachromosomal replication, integration into the host chromosomes or transient expression. Suitable backbones for the present expression vectors include retroviral backbones, backbones from eukaryotic episomes such as SV40 or Bovine Papilloma Virus, backbones from bacterial episomes, or artificial chromosomes.

Preferably, the expression vectors also include a polyA signal downstream of the multiple restriction sites for directing the polyadenylation of mRNA transcribed from the gene inserted into the expression vector.

Following the identification of promoter sequences, proteins which interact with the promoter may be identified as described in Example 56 below.

#### EXAMPLE 56

##### Identification of Proteins Which Interact with

##### Promoter Sequences, Upstream Regulatory Sequences, or mRNA

Sequences within the promoter region which are likely to bind transcription factors may be identified by homology to known transcription factor binding sites or through conventional mutagenesis or deletion analyses of reporter plasmids containing the promoter sequence. For example, deletions may be made in a reporter plasmid containing the promoter sequence of interest operably linked to an assayable reporter gene. The reporter plasmids carrying various deletions within the promoter region are transfected into an appropriate host cell and the effects of the deletions on expression levels is assessed. Transcription factor binding sites within the regions in which deletions reduce expression levels may be further localized using site directed mutagenesis, linker scanning analysis, or other techniques familiar to those skilled in the art.

Nucleic acids encoding proteins which interact with sequences in the promoter may be identified using one-hybrid systems such as those described in the manual accompanying the Matchmaker One-Hybrid System kit available from Clontech (Catalog No. K1603-1). Briefly, the Matchmaker One-hybrid system is used as follows. The target sequence for which it is desired to identify binding proteins is cloned upstream of a selectable reporter gene and integrated into the yeast genome. Preferably, multiple copies of the target sequences are inserted into the reporter plasmid in tandem. A library comprised of fusions between cDNAs to be evaluated for the ability to bind to the promoter and the activation domain of a yeast transcription factor, such as GAL4, is transformed into the yeast strain containing the integrated reporter sequence. The yeast are plated on selective media to select cells

expressing the selectable marker linked to the promoter sequence. The colonies which grow on the selective media contain genes encoding proteins which bind the target sequence. The inserts in the genes encoding the fusion proteins are further characterized by sequencing. In addition, the inserts may be inserted into expression vectors or *in vitro* transcription vectors. Binding of the polypeptides encoded by the inserts to the promoter DNA may be confirmed by techniques familiar to those skilled in the art, such as gel shift analysis or DNase protection analysis.

### VIII. Use of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids in Gene Therapy

The present invention also comprises the use of EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids in gene therapy strategies, including antisense and triple helix strategies as described in Examples 57 and 58 below. In antisense approaches, nucleic acid sequences complementary to an mRNA are hybridized to the mRNA intracellularly, thereby blocking the expression of the protein encoded by the mRNA. The antisense sequences may prevent gene expression through a variety of mechanisms. For example, the antisense sequences may inhibit the ability of ribosomes to translate the mRNA. Alternatively, the antisense sequences may block transport of the mRNA from the nucleus to the cytoplasm, thereby limiting the amount of mRNA available for translation. Another mechanism through which antisense sequences may inhibit gene expression is by interfering with mRNA splicing. In yet another strategy, the antisense nucleic acid may be incorporated in a ribozyme capable of specifically cleaving the target mRNA.

#### EXAMPLE 57

##### Preparation and Use of Antisense Oligonucleotides

The antisense nucleic acid molecules to be used in gene therapy may be either DNA or RNA sequences. They may comprise a sequence complementary to the sequence of the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids. The antisense nucleic acids should have a length and melting temperature sufficient to permit formation of an intracellular duplex with sufficient stability to inhibit the expression of the mRNA in the duplex. Strategies for designing antisense nucleic acids suitable for use in gene therapy are disclosed in Green *et al.*, *Ann. Rev. Biochem.* 55:569-597 (1986) and Izant and Weintraub, *Cell* 36:1007-1015 (1984).

In some strategies, antisense molecules are obtained from a nucleotide sequence encoding a protein by reversing the orientation of the coding region with respect to a promoter so as to transcribe the opposite strand from that which is normally transcribed in the cell. The antisense molecules may be transcribed using *in vitro* transcription systems such as those which employ T7 or SP6 polymerase to

generate the transcript. Another approach involves transcription of the antisense nucleic acids *in vivo* by operably linking DNA containing the antisense sequence to a promoter in an expression vector.

Alternatively, oligonucleotides which are complementary to the strand normally transcribed in the cell may be synthesized *in vitro*. Thus, the antisense nucleic acids are complementary to the

5 corresponding mRNA and are capable of hybridizing to the mRNA to create a duplex. In some embodiments, the antisense sequences may contain modified sugar phosphate backbones to increase stability and make them less sensitive to RNase activity. Examples of modifications suitable for use in antisense strategies are described by Rossi *et al.*, *Pharmacol. Ther.* 50(2):245-254, (1991).

Various types of antisense oligonucleotides complementary to the sequence of the EST-related  
10 nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may be used. In one preferred embodiment, stable and semi-stable antisense oligonucleotides described in International Application No. PCT WO94/23026 are used. In these molecules, the 3' end or both the 3' and 5' ends are engaged in intramolecular hydrogen bonding between complementary base pairs. These molecules are better able to withstand exonuclease attacks  
15 and exhibit increased stability compared to conventional antisense oligonucleotides.

In another preferred embodiment, the antisense oligodeoxynucleotides against herpes simplex virus types 1 and 2 described in International Application No. WO 95/04141 are used.

In yet another preferred embodiment, the covalently cross-linked antisense oligonucleotides described in International Application No. WO 96/31523 are used. These double- or single-stranded  
20 oligonucleotides comprise one or more, respectively, inter- or intra-oligonucleotide covalent cross-linkages, wherein the linkage consists of an amide bond between a primary amine group of one strand and a carboxyl group of the other strand or of the same strand, respectively, the primary amine group being directly substituted in the 2' position of the strand nucleotide monosaccharide ring, and the carboxyl group being carried by an aliphatic spacer group substituted on a nucleotide or nucleotide  
25 analog of the other strand or the same strand, respectively.

The antisense oligodeoxynucleotides and oligonucleotides disclosed in International Application No. WO 92/18522 may also be used. These molecules are stable to degradation and contain at least one transcription control recognition sequence which binds to control proteins and are effective as decoys therefor. These molecules may contain "hairpin" structures, "dumbbell" structures, "modified  
30 dumbbell" structures, "cross-linked" decoy structures and "loop" structures.

In another preferred embodiment, the cyclic double-stranded oligonucleotides described in European Patent Application No. 0 572 287 A2. These ligated oligonucleotide "dumbbells" contain the binding site for a transcription factor and inhibit expression of the gene under control of the transcription factor by sequestering the factor.

35 Use of the closed antisense oligonucleotides disclosed in International Application No. WO 92/19732 is also contemplated. Because these molecules have no free ends, they are more resistant to

degradation by exonucleases than are conventional oligonucleotides. These oligonucleotides may be multifunctional, interacting with several regions which are not adjacent to the target mRNA.

The appropriate level of antisense nucleic acids required to inhibit gene expression may be determined using *in vitro* expression analysis. The antisense molecule may be introduced into the cells  
5 by diffusion, injection, infection or transfection using procedures known in the art. For example, the antisense nucleic acids can be introduced into the body as a bare or naked oligonucleotide, oligonucleotide encapsulated in lipid, oligonucleotide sequence encapsidated by viral protein, or as an oligonucleotide operably linked to a promoter contained in an expression vector. The expression vector may be any of a variety of expression vectors known in the art, including retroviral or viral vectors,  
10 vectors capable of extrachromosomal replication, or integrating vectors. The vectors may be DNA or RNA.

The antisense molecules are introduced onto cell samples at a number of different concentrations preferably between  $1 \times 10^{-10}$  M to  $1 \times 10^{-4}$  M. Once the minimum concentration that can adequately control gene expression is identified, the optimized dose is translated into a dosage suitable  
15 for use *in vivo*. For example, an inhibiting concentration in culture of  $1 \times 10^{-7}$  translates into a dose of approximately 0.6 mg/kg bodyweight. Levels of oligonucleotide approaching 100 mg/kg bodyweight or higher may be possible after testing the toxicity of the oligonucleotide in laboratory animals. It is additionally contemplated that cells from the vertebrate are removed, treated with the antisense oligonucleotide, and reintroduced into the vertebrate.

20 It is further contemplated that the antisense oligonucleotide sequence is incorporated into a ribozyme sequence to enable the antisense to specifically bind and cleave its target mRNA. For technical applications of ribozyme and antisense oligonucleotides see Rossi *et al.*, *supra*.

In a preferred application of this invention, the polypeptide encoded by the gene is first identified, so that the effectiveness of antisense inhibition on translation can be monitored using  
25 techniques that include but are not limited to antibody-mediated tests such as RIAs and ELISA, functional assays, or radiolabeling.

The EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may also be used in gene therapy approaches based on intracellular triple helix formation. Triple helix oligonucleotides are used to inhibit transcription from a  
30 genome. They are particularly useful for studying alterations in cell activity as it is associated with a particular gene. The EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids of the present invention or, more preferably, a portion of those sequences, can be used to inhibit gene expression in individuals having diseases associated with expression of a particular gene. Similarly, the EST-related nucleic acids,  
35 positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids can be used to study the effect of inhibiting transcription of a particular gene within a cell. Traditionally, homopurine sequences were considered the most useful for triple helix strategies.

However, homopyrimidine sequences can also inhibit gene expression. Such homopyrimidine oligonucleotides bind to the major groove at homopurine:homopyrimidine sequences. Thus, both types of sequences from the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids are contemplated within the scope of this invention.

### EXAMPLE 58

#### Preparation and use of Triple Helix Probes

The sequences of the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids are scanned to identify 10-mer to 20-mer homopyrimidine or homopurine stretches which could be used in triple-helix based strategies for inhibiting gene expression. Following identification of candidate homopyrimidine or homopurine stretches, their efficiency in inhibiting gene expression is assessed by introducing varying amounts of oligonucleotides containing the candidate sequences into tissue culture cells which normally express the target gene. The oligonucleotides may be prepared on an oligonucleotide synthesizer or they may be purchased commercially from a company specializing in custom oligonucleotide synthesis, such as GENSET, Paris, France.

The oligonucleotides may be introduced into the cells using a variety of methods known to those skilled in the art, including but not limited to calcium phosphate precipitation, DEAE-Dextran, electroporation, liposome-mediated transfection or native uptake.

Treated cells are monitored for altered cell function or reduced gene expression using techniques such as Northern blotting, RNase protection assays, or PCR based strategies to monitor the transcription levels of the target gene in cells which have been treated with the oligonucleotide. The cell functions to be monitored are predicted based upon the homologies of the target genes corresponding to the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids from which the oligonucleotide were derived with known gene sequences that have been associated with a particular function. The cell functions can also be predicted based on the presence of abnormal physiologies within cells derived from individuals with a particular inherited disease, particularly when the EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids are associated with the disease using techniques described herein.

The oligonucleotides which are effective in inhibiting gene expression in tissue culture cells may then be introduced *in vivo* using the techniques described above and in Example 56 at a dosage calculated based on the *in vitro* results, as described in Example 57.

In some embodiments, the natural (beta) anomers of the oligonucleotide units can be replaced with alpha anomers to render the oligonucleotide more resistant to nucleases. Further, an intercalating agent such as ethidium bromide, or the like, can be attached to the 3' end of the alpha oligonucleotide to

stabilize the triple helix. For information on the generation of oligonucleotides suitable for triple helix formation see Griffin *et al.* (*Science* 245:967-971 (1989)).

#### EXAMPLE 59

5                    Use of EST-related nucleic acids, positional segments of  
                     EST-related nucleic acids or fragments of positional segments of  
                     EST-related nucleic acids to express an Encoded Protein in a Host Organism

                     The EST-related nucleic acids, positional segments of EST-related nucleic acids or fragments of positional segments of EST-related nucleic acids may also be used to express an encoded protein or  
10 polypeptide in a host organism to produce a beneficial effect. In addition, nucleic acids encoding the EST-related polypeptides, positional segments of EST-related polypeptides or fragments of positional segments of EST-related polypeptides may be used to express the encoded protein or polypeptide in a host organism to produce a beneficial effect.

                     In such procedures, the encoded protein or polypeptide may be transiently expressed in the host  
15 organism or stably expressed in the host organism. The encoded protein or polypeptide may have any of the activities described above. The encoded protein or polypeptide may be a protein or polypeptide which the host organism lacks or, alternatively, the encoded protein may augment the existing levels of the protein in the host organism.

                     In some embodiments in which the protein or polypeptide is secreted, nucleic acids encoding the  
20 full length protein (*i.e.* the signal peptide and the mature protein), or nucleic acids encoding only the mature protein (*i.e.* the protein generated when the signal peptide is cleaved off) is introduced into the host organism.

                     The nucleic acids encoding the proteins or polypeptides may be introduced into the host organism using a variety of techniques known to those of skill in the art. For example, the extended  
25 cDNA may be injected into the host organism as naked DNA such that the encoded protein is expressed in the host organism, thereby producing a beneficial effect.

                     Alternatively, the nucleic acids encoding the protein or polypeptide may be cloned into an expression vector downstream of a promoter which is active in the host organism. The expression vector may be any of the expression vectors designed for use in gene therapy, including viral or retroviral  
30 vectors. The expression vector may be directly introduced into the host organism such that the encoded protein is expressed in the host organism to produce a beneficial effect. In another approach, the expression vector may be introduced into cells *in vitro*. Cells containing the expression vector are thereafter selected and introduced into the host organism, where they express the encoded protein or polypeptide to produce a beneficial effect.

35

#### EXAMPLE 60

Use of Signal Peptides To Import Proteins Into Cells

The short core hydrophobic region (h) of signal peptides encoded by the sequences of SEQ ID NOs. 24-728 and 766-792 may also be used as a carrier to import a peptide or a protein of interest, so-called cargo, into tissue culture cells (Lin *et al.*, *J. Biol. Chem.*, 270: 14225-14258 (1995); Du *et al.*, *J. Peptide Res.*, 51: 235-243 (1998); Rojas *et al.*, *Nature Biotech.*, 16: 370-375 (1998)).

5        When cell permeable peptides of limited size (approximately up to 25 amino acids) are to be translocated across cell membrane, chemical synthesis may be used in order to add the h region to either the C-terminus or the N-terminus to the cargo peptide of interest. Alternatively, when longer peptides or proteins are to be imported into cells, nucleic acids can be genetically engineered, using techniques familiar to those skilled in the art, in order to link the extended cDNA sequence encoding the h region to  
10    the 5' or the 3' end of a DNA sequence coding for a cargo polypeptide. Such genetically engineered nucleic acids are then translated either *in vitro* or *in vivo* after transfection into appropriate cells, using conventional techniques to produce the resulting cell permeable polypeptide. Suitable hosts cells are then simply incubated with the cell permeable polypeptide which is then translocated across the membrane.

15        This method may be applied to study diverse intracellular functions and cellular processes. For instance, it has been used to probe functionally relevant domains of intracellular proteins and to examine protein-protein interactions involved in signal transduction pathways (Lin *et al.*, *supra*; Lin *et al.*, *J. Biol. Chem.*, 271: 5305-5308 (1996); Rojas *et al.*, *J. Biol. Chem.*, 271: 27456-27461 (1996); Liu *et al.*, *Proc. Natl. Acad. Sci. USA*, 93: 11819-11824 (1996); Rojas *et al.*, *Bioch. Biophys. Res. Commun.*, 234: 675-  
20    680 (1997)).

Such techniques may be used in cellular therapy to import proteins producing therapeutic effects. For instance, cells isolated from a patient may be treated with imported therapeutic proteins and then re-introduced into the host organism.

Alternatively, the h region of signal peptides of the present invention could be used in  
25    combination with a nuclear localization signal to deliver nucleic acids into cell nucleus. Such oligonucleotides may be antisense oligonucleotides or oligonucleotides designed to form triple helixes, as described above, in order to inhibit processing and maturation of a target cellular RNA.

### EXAMPLE 61

#### 30        Computer Embodiments

As used herein the term "nucleic acid codes of SEQ ID NOs. 24-811 and 1600-1622" encompasses the nucleotide sequences of SEQ ID NOs. 24-811 and 1600-1622, fragments of SEQ ID NOs. 24-811 and 1600-1622, nucleotide sequences homologous to SEQ ID NOs. 24-811 and 1600-1622 or homologous to fragments of SEQ ID NOs. 24-811 and 1600-1622, and sequences  
35    complementary to all of the preceding sequences. The fragments include portions of SEQ ID NOs. 24-811 and 1600-1622 comprising at least 10, 15, 20, 25, 30, 35, 40, 50, 75, 100, 150, 200, 300, 400, or 500 consecutive nucleotides of SEQ ID NOs. 24-811 and 1600-1622. Preferably, the fragments are novel



fragments. Preferably the fragments include polynucleotides described in Table II, polynucleotides described in Table III, polynucleotides described in Table IV or portions thereof comprising at least 10, 15, 20, 25, 30, 35, 40, 50, 75, 100, 150, 200, 300, 400, or 500 consecutive nucleotides of the polynucleotides described in Tables II, III, or IV. Homologous sequences and fragments of SEQ ID NOs. 24-811 and 1600-1622 refer to a sequence having at least 99%, 98%, 97%, 96%, 95%, 90%, 85%, 80%, or 75% homology to these sequences. Homology may be determined using any of the computer programs and parameters described in Example 18, including BLAST2N with the default parameters or with any modified parameters. Homologous sequences also include RNA sequences in which uridines replace the thymines in the nucleic acid codes of SEQ ID NOs. 24-811 and 1600-1622. The homologous sequences may be obtained using any of the procedures described herein or may result from the correction of a sequencing error as described above. Preferably the homologous sequences and fragments of SEQ ID NOs. 24-811 and 1600-1622 include polynucleotides described in Table II, polynucleotides described in Table III, polynucleotides described in Table IV or portions thereof comprising at least 10, 15, 20, 25, 30, 35, 40, 50, 75, 100, 150, 200, 300, 400, or 500 consecutive nucleotides of the polynucleotides described in Tables II, III, or IV. It will be appreciated that the nucleic acid codes of SEQ ID NOs. 24-811 and 1600-1622 can be represented in the traditional single character format (See the inside back cover of *Styer, Lubert. Biochemistry*, 3<sup>rd</sup> edition. W. H Freeman & Co., New York.) or in any other format which records the identity of the nucleotides in a sequence.

As used herein the term "polypeptide codes of SEQ ID NOS. 812-1599" encompasses the polypeptide sequence of SEQ ID NOS. 812-1599 which are encoded by the 5' ESTs of SEQ ID NOS. 24-811 and 1600-1622, polypeptide sequences homologous to the polypeptides of SEQ ID NOS. 812-1599, or fragments of any of the preceding sequences. Homologous polypeptide sequences refer to a polypeptide sequence having at least 99%, 98%, 97%, 96%, 95%, 90%, 85%, 80%, 75% homology to one of the polypeptide sequences of SEQ ID NOS. 812-1599. Homology may be determined using any of the computer programs and parameters described herein, including FASTA with the default parameters or with any modified parameters. The homologous sequences may be obtained using any of the procedures described herein or may result from the correction of a sequencing error as described above. The polypeptide fragments comprise at least 5, 10, 15, 20, 25, 30, 35, 40, 50, 75, 100, or 150 consecutive amino acids of the polypeptides of SEQ ID NOS. 812-1599. Preferably, the fragments are novel fragments. Preferably, the fragments include polypeptides encoded by the polynucleotides described in Table II, or portions thereof comprising at least 5, 10, 15, 20, 25, 30, 35, 40, 50, 75, 100, or 150 consecutive amino acids of the polypeptides encoded by the polynucleotides described in Table II. It will be appreciated that the polypeptide codes of the SEQ ID NOS. 812-1599 can be represented in the traditional single character format or three letter format (See the inside back cover of *Starrier, Lubert. Biochemistry*, 3<sup>rd</sup> edition. W. H Freeman & Co., New York.) or in any other format which relates the identity of the polypeptides in a sequence.

It will be appreciated by those skilled in the art that the nucleic acid codes of SEQ ID NOS. 24-811 and 1600-1622 and polypeptide codes of SEQ ID NOS. 812-1599 can be stored, recorded, and manipulated on any medium which can be read and accessed by a computer. As used herein, the words "recorded" and "stored" refer to a process for storing information on a computer medium. A skilled artisan can readily adopt any of the presently known methods for recording information on a computer readable medium to generate manufactures comprising one or more of the nucleic acid codes of SEQ ID NOS. 24-811 and 1600-1622, one or more of the polypeptide codes of SEQ ID NOS. 812-1599. Another aspect of the present invention is a computer readable medium having recorded thereon at least 2, 5, 10, 15, 20, 25, 30, or 50 nucleic acid codes of SEQ ID NOS. 24-811 and 1600-1622. Another aspect of the present invention is a computer readable medium having recorded thereon at least 2, 5, 10, 15, 20, 25, 30, or 50 polypeptide codes of SEQ ID NOS. 812-1599.

Computer readable media include magnetically readable media, optically readable media, electronically readable media and magnetic/optical media. For example, the computer readable media may be a hard disk, a floppy disk, a magnetic tape, CD-ROM, Digital Versatile Disk (DVD), Random Access Memory (RAM), or Read Only Memory (ROM) as well as other types of other media known to those skilled in the art.

Embodiments of the present invention include systems, particularly computer systems which store and manipulate the sequence information described herein. One example of a computer system 100 is illustrated in block diagram form in Figure 6. As used herein, "a computer system" refers to the hardware components, software components, and data storage components used to analyze the nucleotide sequences of the nucleic acid codes of SEQ ID NOS. 24-811 and 1600-1622, or the amino acid sequences of the polypeptide codes of SEQ ID NOS. 812-1599. In one embodiment, the computer system 100 is a Sun Enterprise 1000 server (Sun Microsystems, Palo Alto, CA). The computer system 100 preferably includes a processor for processing, accessing and manipulating the sequence data. The processor 105 can be any well-known type of central processing unit, such as the Pentium III from Intel Corporation, or similar processor from Sun, Motorola, Compaq or International Business Machines.

Preferably, the computer system 100 is a general purpose system that comprises the processor 105 and one or more internal data storage components 110 for storing data, and one or more data retrieving devices for retrieving the data stored on the data storage components. A skilled artisan can readily appreciate that any one of the currently available computer systems are suitable.

In one particular embodiment, the computer system 100 includes a processor 105 connected to a bus which is connected to a main memory 115 (preferably implemented as RAM) and one or more internal data storage devices 110, such as a hard drive and/or other computer readable media having data recorded thereon. In some embodiments, the computer system 100 further includes one or more data retrieving device 118 for reading the data stored on the internal data storage devices 110.

The data retrieving device 118 may represent, for example, a floppy disk drive, a compact disk drive, a magnetic tape drive, etc. In some embodiments, the internal data storage device 110 is a

removable computer readable medium such as a floppy disk, a compact disk, a magnetic tape, etc. containing control logic and/or data recorded thereon. The computer system 100 may advantageously include or be programmed by appropriate software for reading the control logic and/or the data from the data storage component once inserted in the data retrieving device.

5       The computer system 100 includes a display 120 which is used to display output to a computer user. It should also be noted that the computer system 100 can be linked to other computer systems 125a-c in a network or wide area network to provide centralized access to the computer system 100.

Software for accessing and processing the nucleotide sequences of the nucleic acid codes of SEQ ID NOS. 24-811 and 1600-1622, or the amino acid sequences of the polypeptide codes of SEQ ID  
10 NOS. 812-1599 (such as search tools, compare tools, and modeling tools etc.) may reside in main memory 115 during execution.

In some embodiments, the computer system 100 may further comprise a sequence comparer for comparing the above-described nucleic acid codes of SEQ ID NOS. 24-811 and 1600-1622 or polypeptide codes of SEQ ID NOS. 812-1599 stored on a computer readable medium to reference  
15 nucleotide or polypeptide sequences stored on a computer readable medium. A "sequence comparer" refers to one or more programs which are implemented on the computer system 100 to compare a nucleotide or polypeptide sequence with other nucleotide or polypeptide sequences and/or compounds including but not limited to peptides, peptidomimetics, and chemicals stored within the data storage means. For example, the sequence comparer may compare the nucleotide sequences of the nucleic acid  
20 codes of SEQ ID NOS. 24-811 and 1600-1622, or the amino acid sequences of the polypeptide codes of SEQ ID NOS. 812-1599 stored on a computer readable medium to reference sequences stored on a computer readable medium to identify homologies, motifs implicated in biological function, or structural motifs. The various sequence comparer programs identified elsewhere in this patent specification are particularly contemplated for use in this aspect of the invention.

25       Figure 7 is a flow diagram illustrating one embodiment of a process 200 for comparing a new nucleotide or protein sequence with a database of sequences in order to determine the homology levels between the new sequence and the sequences in the database. The database of sequences can be a private database stored within the computer system 100, or a public database such as GENBANK, PIR OR SWISSPROT that is available through the Internet.

30       The process 200 begins at a start state 201 and then moves to a state 202 wherein the new sequence to be compared is stored to a memory in a computer system 100. As discussed above, the memory could be any type of memory, including RAM or an internal storage device.

The process 200 then moves to a state 204 wherein a database of sequences is opened for analysis and comparison. The process 200 then moves to a state 206 wherein the first sequence stored in  
35 the database is read into a memory on the computer. A comparison is then performed at a state 210 to determine if the first sequence is the same as the second sequence. It is important to note that this step is not limited to performing an exact comparison between the new sequence and the first sequence in the

database. Well-known methods are known to those of skill in the art for comparing two nucleotide or protein sequences, even if they are not identical. For example, gaps can be introduced into one sequence in order to raise the homology level between the two tested sequences. The parameters that control whether gaps or other features are introduced into a sequence during comparison are normally entered by the user of the computer system.

Once a comparison of the two sequences has been performed at the state 210, a determination is made at a decision state 210 whether the two sequences are the same. Of course, the term "same" is not limited to sequences that are absolutely identical. Sequences that are within the homology parameters entered by the user will be marked as "same" in the process 200.

If a determination is made that the two sequences are the same, the process 200 moves to a state 214 wherein the name of the sequence from the database is displayed to the user. This state notifies the user that the sequence with the displayed name fulfills the homology constraints that were entered. Once the name of the stored sequence is displayed to the user, the process 200 moves to a decision state 218 wherein a determination is made whether more sequences exist in the database. If no more sequences exist in the database, then the process 200 terminates at an end state 220. However, if more sequences do exist in the database, then the process 200 moves to a state 224 wherein a pointer is moved to the next sequence in the database so that it can be compared to the new sequence. In this manner, the new sequence is aligned and compared with every sequence in the database.

It should be noted that if a determination had been made at the decision state 212 that the sequences were not homologous, then the process 200 would move immediately to the decision state 218 in order to determine if any other sequences were available in the database for comparison.

Accordingly, one aspect of the present invention is a computer system comprising a processor, a data storage device having stored thereon a nucleic acid code of SEQ ID NOS. 24-811 and 1600-1622 or a polypeptide code of SEQ ID NOS. 812-1599, a data storage device having retrievably stored thereon reference nucleotide sequences or polypeptide sequences to be compared to the nucleic acid code of SEQ ID NOS. 24-811 and 1600-1622 or polypeptide code of SEQ ID NOS. 812-1599 and a sequence comparer for conducting the comparison. The sequence comparer may indicate a homology level between the sequences compared or identify structural motifs in the above described nucleic acid code of SEQ ID NOS. 24-811 and 1600-1622 and polypeptide codes of SEQ ID NOS. 812-1599 or it may identify structural motifs in sequences which are compared to these nucleic acid codes and polypeptide codes. In some embodiments, the data storage device may have stored thereon the sequences of at least 2, 5, 10, 15, 20, 25, 30, or 50 of the nucleic acid codes of SEQ ID NOS. 24-811 and 1600-1622 or polypeptide codes of SEQ ID NOS. 812-1599.

Another aspect of the present invention is a method for determining the level of homology between a nucleic acid code of SEQ ID NOS. 24-811 and 1600-1622 and a reference nucleotide sequence, comprising the steps of reading the nucleic acid code and the reference nucleotide sequence through the use of a computer program which determines homology levels and determining homology

between the nucleic acid code and the reference nucleotide sequence with the computer program. The computer program may be any of a number of computer programs for determining homology levels, including those specifically enumerated herein, including BLAST2N with the default parameters or with any modified parameters. The method may be implemented using the computer systems described above. The method may also be performed by reading 2, 5, 10, 15, 20, 25, 30, or 50 of the above described nucleic acid codes of SEQ ID NOs. 24-811 and 1600-1622 through use of the computer program and determining homology between the nucleic acid codes and reference nucleotide sequences.

Figure 8 is a flow diagram illustrating one embodiment of a process 250 in a computer for determining whether two sequences are homologous. The process 250 begins at a start state 252 and then moves to a state 254 wherein a first sequence to be compared is stored to a memory. The second sequence to be compared is then stored to a memory at a state 256. The process 250 then moves to a state 260 wherein the first character in the first sequence is read and then to a state 262 wherein the first character of the second sequence is read. It should be understood that if the sequence is a nucleotide sequence, then the character would normally be either A, T, C, G or U. If the sequence is a protein sequence, then it should be in the single letter amino acid code so that the first and sequence sequences can be easily compared.

A determination is then made at a decision state 264 whether the two characters are the same. If they are the same, then the process 250 moves to a state 268 wherein the next characters in the first and second sequences are read. A determination is then made whether the next characters are the same. If they are, then the process 250 continues this loop until two characters are not the same. If a determination is made that the next two characters are not the same, the process 250 moves to a decision state 274 to determine whether there are any more characters either sequence to read.

If there aren't any more characters to read, then the process 250 moves to a state 276 wherein the level of homology between the first and second sequences is displayed to the user. The level of homology is determined by calculating the proportion of characters between the sequences that were the same out of the total number of sequences in the first sequence. Thus, if every character in a first 100 nucleotide sequence aligned with a every character in a second sequence, the homology level would be 100%.

Alternatively, the computer program may be a computer program which compares the nucleotide sequences of the nucleic acid codes of the present invention, to reference nucleotide sequences in order to determine whether the nucleic acid code of SEQ ID NOs. 24-811 and 1600-1622 differs from a reference nucleic acid sequence at one or more positions. Optionally such a program records the length and identity of inserted, deleted or substituted nucleotides with respect to the sequence of either the reference polynucleotide or the nucleic acid code of SEQ ID NOs. 24-811 and 1600-1622. In one embodiment, the computer program may be a program which determines whether the nucleotide sequences of the nucleic acid codes of SEQ ID NOs. 24-811 and 1600-1622 contain a biallelic marker

or single nucleotide polymorphism (SNP) with respect to a reference nucleotide sequence. This single nucleotide polymorphism may comprise a single base substitution, insertion, or deletion, while this biallelic marker may comprise about one to ten consecutive bases substituted, inserted or deleted.

Another aspect of the present invention is a method for determining the level of homology  
5 between a polypeptide code of SEQ ID NOS. 812-1599 and a reference polypeptide sequence, comprising the steps of reading the polypeptide code of SEQ ID NOS. 812-1599 and the reference polypeptide sequence through use of a computer program which determines homology levels and determining homology between the polypeptide code and the reference polypeptide sequence using the computer program.

10 Accordingly, another aspect of the present invention is a method for determining whether a nucleic acid code of SEQ ID NOS. 24-811 and 1600-1622 differs at one or more nucleotides from a reference nucleotide sequence comprising the steps of reading the nucleic acid code and the reference nucleotide sequence through use of a computer program which identifies differences between nucleic acid sequences and identifying differences between the nucleic acid code and the reference nucleotide  
15 sequence with the computer program. In some embodiments, the computer program is a program which identifies single nucleotide polymorphisms. The method may be implemented by the computer systems described above and the method illustrated in Figure 8. The method may also be performed by reading at least 2, 5, 10, 15, 20, 25, 30, or 50 of the nucleic acid codes of SEQ ID NOS. 24-811 and 1600-1622 and the reference nucleotide sequences through the use of the computer program and identifying differences  
20 between the nucleic acid codes and the reference nucleotide sequences with the computer program.

In other embodiments the computer based system may further comprise an identifier for identifying features within the nucleotide sequences of the nucleic acid codes of SEQ ID NOS. 24-811 and 1600-1622 or the amino acid sequences of the polypeptide codes of SEQ ID NOS. 812-1599.

An "identifier" refers to one or more programs which identifies certain features within the  
25 above-described nucleotide sequences of the nucleic acid codes of SEQ ID NOS. 24-811 and 1600-1622 or the amino acid sequences of the polypeptide codes of SEQ ID NOS. 812-1599. In one embodiment, the identifier may comprise a program which identifies an open reading frame in the cDNAs codes of SEQ ID NOS. 24-811 and 1600-1622.

Figure 9 is a flow diagram illustrating one embodiment of an identifier process 300 for  
30 detecting the presence of a feature in a sequence. The process 300 begins at a start state 302 and then moves to a state 304 wherein a first sequence that is to be checked for features is stored to a memory 115 in the computer system 100. The process 300 then moves to a state 306 wherein a database of sequence features is opened. Such a database would include a list of each feature's attributes along with the name of the feature. For example, a feature name could be "Initiation Codon" and the  
35 attribute would be "ATG". Another example would be the feature name "TAATAA Box" and the feature attribute would be "TAATAA". An example of such a database is produced by the University of Wisconsin Genetics Computer Group ([www.gcg.com](http://www.gcg.com)).

Once the database of features is opened at the state 306, the process 300 moves to a state 308 wherein the first feature is read from the database. A comparison of the attribute of the first feature with the first sequence is then made at a state 310. A determination is then made at a decision state 316 whether the attribute of the feature was found in the first sequence. If the attribute was found, then the process 300 moves to a state 318 wherein the name of the found feature is displayed to the user.

The process 300 then moves to a decision state 320 wherein a determination is made whether more features exist in the database. If no more features do exist, then the process 300 terminates at an end state 324. However, if more features do exist in the database, then the process 300 reads the next sequence feature at a state 326 and loops back to the state 310 wherein the attribute of the next feature is compared against the first sequence.

It should be noted, that if the feature attribute is not found in the first sequence at the decision state 316, the process 300 moves directly to the decision state 320 in order to determine if any more features exist in the database.

In another embodiment, the identifier may comprise a molecular modeling program which determines the 3-dimensional structure of the polypeptides codes of SEQ ID NOS. 812-1599. In some embodiments, the molecular modeling program identifies target sequences that are most compatible with profiles representing the structural environments of the residues in known three-dimensional protein structures. (See, e.g., Eisenberg et al., U.S. Patent No. 5,436,850 issued July 25, 1995). In another technique, the known three-dimensional structures of proteins in a given family are superimposed to define the structurally conserved regions in that family. This protein modeling technique also uses the known three-dimensional structure of a homologous protein to approximate the structure of the polypeptide codes of SEQ ID NOS. 812-1599. (See e.g., Srinivasan, et al., U.S. Patent No. 5,557,535 issued September 17, 1996). Conventional homology modeling techniques have been used routinely to build models of proteases and antibodies. (Sowdhamini et al., Protein Engineering 10:207, 215 (1997)). Comparative approaches can also be used to develop three-dimensional protein models when the protein of interest has poor sequence identity to template proteins. In some cases, proteins fold into similar three-dimensional structures despite having very weak sequence identities. For example, the three-dimensional structures of a number of helical cytokines fold in similar three-dimensional topology in spite of weak sequence homology.

The recent development of threading methods now enables the identification of likely folding patterns in a number of situations where the structural relatedness between target and template(s) is not detectable at the sequence level. Hybrid methods, in which fold recognition is performed using Multiple Sequence Threading (MST), structural equivalencies are deduced from the threading output using a distance geometry program DRAGON to construct a low resolution model, and a full-atom representation is constructed using a molecular modeling package such as QUANTA.

According to this 3-step approach, candidate templates are first identified by using the novel fold recognition algorithm MST, which is capable of performing simultaneous threading of multiple aligned sequences onto one or more 3-D structures. In a second step, the structural equivalencies obtained from the MST output are converted into interresidue distance restraints and fed into the distance geometry program DRAGON, together with auxiliary information obtained from secondary structure predictions. The program combines the restraints in an unbiased manner and rapidly generates a large number of low resolution model confirmations. In a third step, these low resolution model confirmations are converted into full-atom models and subjected to energy minimization using the molecular modeling package QUANTA. (See e.g., Aszódi et al., *Proteins: Structure, Function, and Genetics*, Supplement 1:38-42 (1997)).

The results of the molecular modeling analysis may then be used in rational drug design techniques to identify agents which modulate the activity of the polypeptide codes of SEQ ID NOS. 812-1599.

Accordingly, another aspect of the present invention is a method of identifying a feature within the nucleic acid codes of SEQ ID NOS. 24-811 and 1600-1622 or the polypeptide codes of SEQ ID NOS. 812-1599 comprising reading the nucleic acid code(s) or the polypeptide code(s) through the use of a computer program which identifies features therein and identifying features within the nucleic acid code(s) or polypeptide code(s) with the computer program. In one embodiment, computer program comprises a computer program which identifies open reading frames. In a further embodiment, the computer program identifies structural motifs in a polypeptide sequence. In another embodiment, the computer program comprises a molecular modeling program. The method may be performed by reading a single sequence or at least 2, 5, 10, 15, 20, 25, 30, or 50 of the nucleic acid codes of SEQ ID NOS. 24-811 and 1600-1622 or the polypeptide codes of SEQ ID NOS. 812-1599 through the use of the computer program and identifying features within the nucleic acid codes or polypeptide codes with the computer program.

The nucleic acid codes of SEQ ID NOS. 24-811 and 1600-1622 or the polypeptide codes of SEQ ID NOS. 812-1599 may be stored and manipulated in a variety of data processor programs in a variety of formats. For example, the nucleic acid codes of SEQ ID NOS. 24-811 and 1600-1622 or the polypeptide codes of SEQ ID NOS. 812-1599 may be stored as text in a word processing file, such as MicrosoftWORD or WORDPERFECT or as an ASCII file in a variety of database programs familiar to those of skill in the art, such as DB2, SYBASE, or ORACLE. In addition, many computer programs and databases may be used as sequence comparers, identifiers, or sources of reference nucleotide or polypeptide sequences to be compared to the nucleic acid codes of SEQ ID NOS. 24-811 and 1600-1622 or the polypeptide codes of SEQ ID NOS. 812-1599. The following list is intended not to limit the invention but to provide guidance to programs and databases which are useful with the nucleic acid codes of SEQ ID NOS. 24-811 and 1600-1622 or the polypeptide codes of SEQ ID NOS. 812-1599. The programs and databases which may be used include, but are not limited to: MacPattern (EMBL),



DiscoveryBase (Molecular Applications Group), GeneMine (Molecular Applications Group), Look (Molecular Applications Group), MacLook (Molecular Applications Group), BLAST and BLAST2 (NCBI), BLASTN and BLASTX (Altschul et al, *J. Mol. Biol.* 215: 403 (1990)), FASTA (Pearson and Lipman, *Proc. Natl. Acad. Sci. USA*, 85: 2444 (1988)), FASTDB (Brutlag et al. *Comp. App. Biosci.* 6:237-245, 1990), Catalyst (Molecular Simulations Inc.), Catalyst/SHAPE (Molecular Simulations Inc.), Cerius<sup>2</sup>.DBAccess (Molecular Simulations Inc.), HypoGen (Molecular Simulations Inc.), Insight II, (Molecular Simulations Inc.), Discover (Molecular Simulations Inc.), CHARMM (Molecular Simulations Inc.), Felix (Molecular Simulations Inc.), DelPhi, (Molecular Simulations Inc.), QuanteMM, (Molecular Simulations Inc.), Homology (Molecular Simulations Inc.), Modeler (Molecular Simulations Inc.), ISIS (Molecular Simulations Inc.), Quanta/Protein Design (Molecular Simulations Inc.), WebLab (Molecular Simulations Inc.), WebLab Diversity Explorer (Molecular Simulations Inc.), Gene Explorer (Molecular Simulations Inc.), SeqFold (Molecular Simulations Inc.), the EMBL/Swissprotein database, the MDL Available Chemicals Directory database, the MDL Drug Data Report data base, the Comprehensive Medicinal Chemistry database, Derwents's World Drug Index database, the BioByteMasterFile database, the Genbank database, and the Genseqn database. Many other programs and data bases would be apparent to one of skill in the art given the present disclosure.

Motifs which may be detected using the above programs include sequences encoding leucine zippers, helix-turn-helix motifs, glycosylation sites, ubiquitination sites, alpha helices, and beta sheets, signal sequences encoding signal peptides which direct the secretion of the encoded proteins, sequences implicated in transcription regulation such as homeoboxes, acidic stretches, enzymatic active sites, substrate binding sites, and enzymatic cleavage sites.

## EXAMPLE 62

### Methods of Making Nucleic Acids

The present invention also comprises methods of making the EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of the EST-related nucleic acids, or fragments of positional segments of the EST-related nucleic acids. The methods comprise sequentially linking together nucleotides to produce the nucleic acids having the preceding sequences. A variety of methods of synthesizing nucleic acids are known to those skilled in the art.

In many of these methods, synthesis is conducted on a solid support. These included the 3' phosphoramidite methods in which the 3' terminal base of the desired oligonucleotide is immobilized on an insoluble carrier. The nucleotide base to be added is blocked at the 5' hydroxyl and activated at the 3' hydroxyl so as to cause coupling with the immobilized nucleotide base. Deblocking of the new immobilized nucleotide compound and repetition of the cycle will produce the desired polynucleotide. Alternatively, polynucleotides may be prepared as described in U.S. Patent No. 5,049,656. In some embodiments, several polynucleotides prepared as described above are ligated together to generate longer polynucleotides having a desired sequence.

**EXAMPLE 63****Methods of Making Polypeptides**

The present invention also comprises methods of making the polynucleotides encoded by EST-related nucleic acids, fragments of EST-related nucleic acids, positional segments of the EST-related nucleic acids, or fragments of positional segments of the EST-related nucleic acids and methods of making the EST-related polypeptides, fragments of EST-related polypeptides, positional segments of EST-related polypeptides, or fragments of EST-related polypeptides. The methods comprise sequentially linking together amino acids to produce the nucleic polypeptides having the preceding sequences. In some embodiments, the polypeptides made by these methods are 150 amino acid or less in length. In other embodiments, the polypeptides made by these methods are 120 amino acids or less in length.

A variety of methods of making polypeptides are known to those skilled in the art, including methods in which the carboxyl terminal amino acid is bound to polyvinyl benzene or another suitable resin. The amino acid to be added possesses blocking groups on its amino moiety and any side chain reactive groups so that only its carboxyl moiety can react. The carboxyl group is activated with carbodiimide or another activating agent and allowed to couple to the immobilized amino acid. After removal of the blocking group, the cycle is repeated to generate a polypeptide having the desired sequence. Alternatively, the methods described in U.S. Patent No. 5,049,656 may be used.

As discussed above, the EST-related nucleic acids, fragments of the EST-related nucleic acids, positional segments of the EST-related nucleic acids, or fragments of positional segments of the EST-related nucleic acids can be used for various purposes. The polynucleotides can be used to express recombinant protein for analysis, characterization or therapeutic use; production of secreted polypeptides or chimeric polypeptides, antibody production, as markers for tissues in which the corresponding protein is preferentially expressed (either constitutively or at a particular stage of tissue differentiation or development or in disease states); as molecular weight markers on Southern gels; as chromosome markers or tags (when labeled) to identify chromosomes or to map related gene positions; to compare with endogenous DNA sequences in patients to identify potential genetic disorders; as probes to hybridize and thus discover novel, related DNA sequences; as a source of information to derive PCR primers for genetic fingerprinting; for selecting and making oligomers for attachment to a "gene chip" or other support, including for examination for expression patterns; to raise anti-protein antibodies using DNA immunization techniques; and as an antigen to raise anti-DNA antibodies or elicit another immune response. Where the polynucleotide encodes a protein or polypeptide which binds or potentially binds to another protein or polypeptide (such as, for example, in a receptor-ligand interaction), the polynucleotide can also be used in interaction trap assays (such as, for example, that described in Gyuris *et al.*, *Cell* 75:791-803 (1993)) to identify polynucleotides encoding the other protein or polypeptide with which binding occurs or to identify inhibitors of the binding interaction.

The proteins or polypeptides provided by the present invention can similarly be used in assays to determine biological activity, including in a panel of multiple proteins for high-throughput screening; to raise antibodies or to elicit another immune response; as a reagent (including the labeled reagent) in assays designed to quantitatively determine levels of the protein (or its receptor) in biological fluids; as markers for tissues in which the corresponding protein is preferentially expressed (either constitutively or at a particular stage of tissue differentiation or development or in a disease state); and, of course, to isolate correlative receptors or ligands. Where the protein or polypeptide binds or potentially binds to another protein or polypeptide (such as, for example, in a receptor-ligand interaction), the protein can be used to identify the other protein with which binding occurs or to identify inhibitors of the binding interaction. Proteins or polypeptides involved in these binding interactions can also be used to screen for peptide or small molecule inhibitors or agonists of the binding interaction.

Any or all of these research utilities are capable of being developed into reagent grade or kit format for commercialization as research products.

Methods for performing the uses listed above are well known to those skilled in the art.

References disclosing such methods include without limitation "Molecular Cloning; A Laboratory Manual," 2d ed., Cold Spring Harbor Laboratory Press, Sambrook, J., E.F. Fritsch and T. Maniatis eds., 1989, and "Methods in Enzymology; Guide to Molecular Cloning Techniques," Academic Press, Berger, S.L. and A.R. Kimmel eds., 1987.

Polynucleotides and proteins or polypeptides of the present invention can also be used as nutritional sources or supplements. Such uses include without limitation use as a protein or amino acid supplement, use as a carbon source, use as a nitrogen source and use as a source of carbohydrate. In such cases the protein or polynucleotide of the invention can be added to the feed of a particular organism or can be administered as a separate solid or liquid preparation, such as in the form of powder, pills, solutions, suspensions or capsules. In the case of microorganisms, the protein or polynucleotide of the invention can be added to the medium in or on which the microorganism is cultured.

Although this invention has been described in terms of certain preferred embodiments, other embodiments which will be apparent to those of ordinary skill in the art in view of the disclosure herein are also within the scope of this invention. Accordingly, the scope of the invention is intended to be limited only by reference to the appended claims.

#### Sequence Listing Free Text

The following free text appears in the accompanying Sequence Listing:

Von Heijne matrix

score

sequence

name

martinspector prediction

CLAIMS

1. A purified nucleic acid comprising a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 and sequences complementary to the sequences of  
5 SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622.

2. A purified nucleic acid comprising at least 15 consecutive nucleotides of a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 and sequences complementary to the sequences of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622.

10

3. A purified or isolated polypeptide comprising a sequence selected from the group consisting of the sequences of SEQ ID NOs. 812-1599.

4. A method of making a cDNA comprising the steps of:

15

a) contacting a collection of mRNA molecules from human cells with a primer comprising at least 15 consecutive nucleotides of a sequence selected from the group consisting of the sequences complementary to SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622;

b) hybridizing said primer to an mRNA in said collection that encodes said protein;

20

c) reverse transcribing said hybridized primer to make a first cDNA strand from said mRNA;

d) making a second cDNA strand complementary to said first cDNA strand; and

e) isolating the resulting cDNA comprising said first cDNA strand and said second cDNA strand.

25

5. A method of making a cDNA comprising the steps of:

a) obtaining a cDNA comprising a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622;

b) contacting said cDNA with a detectable probe comprising at least 15 consecutive nucleotides of a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID

30

NOs. 1600-1622 and the sequences complementary to SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 under conditions which permit said probe to hybridize to said cDNA;

c) identifying a cDNA which hybridizes to said detectable probe; and

d) isolating said cDNA which hybridizes to said probe.

35

6. A method of making a cDNA comprising the steps of:

a) contacting a collection of mRNA molecules from human cells with a first primer capable of hybridizing to the polyA tail of said mRNA;

b) hybridizing said first primer to said polyA tail;

- c) reverse transcribing said mRNA to make a first cDNA strand;
- d) making a second cDNA strand complementary to said first cDNA strand using at least one primer comprising at least 15 consecutive nucleotides of a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622; and
- 5 e) isolating the resulting cDNA comprising said first cDNA strand and said second cDNA strand.

7. A method of making a polypeptide comprising the steps of:

- a) obtaining a cDNA which encodes a polypeptide encoded by a nucleic acid comprising  
10 a sequence selected from the group consisting of SEQ ID NOs. 24-811 or a cDNA which encodes a polypeptide comprising at least 10 consecutive amino acids of a polypeptide encoded by a sequence selected from the group consisting of SEQ ID NOs. 24-811;
- b) inserting said cDNA in an expression vector such that said cDNA is operably linked to a promoter;
- 15 c) introducing said expression vector into a host cell whereby said host cell produces the protein encoded by said cDNA; and
- d) isolating said protein.

8. In an array of discrete ESTs or fragments thereof of at least 15 nucleotides in length, the  
20 improvement comprising inclusion in said array of at least one sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622, the sequences complementary to the sequences of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 and fragments comprising at least 15 consecutive nucleotides of said sequence.

25 9. The array of Claim 8 including therein at least five sequences selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622, the sequences complementary to the sequences of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 and fragments comprising at least 15 consecutive nucleotides of said sequences.

30 10. An enriched population of recombinant nucleic acids, said recombinant nucleic acids comprising an insert nucleic acid and a backbone nucleic acid, wherein at least 5% of said insert nucleic acids in said population comprise a sequence selected from the group consisting of SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622, the sequences complementary to SEQ ID NOs. 24-811 and SEQ ID NOs. 1600-1622 and fragments comprising at least 15 consecutive nucleotides of said  
35 sequences.

11. An antibody composition capable of selectively binding to an epitope-containing fragment of a polypeptide comprising a contiguous span of at least 8 amino acids of any of SEQ ID NOs. 812-1599, wherein said antibody is polyclonal or monoclonal.

12. A computer readable medium having stored thereon a sequence selected from the group consisting of a nucleic acid code of SEQ ID NOs. 24-811 and 1600-1622 and a polypeptide code of SEQ ID NOs. 812-1599.

5

13. A computer system comprising a processor and a data storage device wherein said data storage device has stored thereon a sequence selected from the group consisting of a nucleic acid code of SEQ ID NOs. 24-811 and 1600-1622 and a polypeptide code of SEQ ID NOs. 812-1599.

10

14. The computer system of Claim 13 further comprising a sequence comparer and a data storage device having reference sequences stored thereon.

15. The computer system of Claim 14 wherein said sequence comparer comprises a computer program which indicates polymorphisms.

15

16. The computer system of Claim 13 further comprising an identifier which identifies features in said sequence.

17. A method for comparing a first sequence to a reference sequence wherein said first sequence is selected from the group consisting of a nucleic acid code of SEQ ID NOs. 24-811 and 1600-1622 and a polypeptide code of SEQ ID NOs. 812-1599 comprising the steps of:

- a) reading said first sequence and said reference sequence through use of a computer program which compares sequences; and
- b) determining differences between said first sequence and said reference sequence with said computer program.

25

18. The method of Claim 17, wherein said step of determining differences between the first sequence and the reference sequence comprises identifying polymorphisms.

30

19. A method for identifying a feature in a sequence selected from the group consisting of a nucleic acid code of SEQ ID NOs. 24-811 and 1600-1622 and a polypeptide code of SEQ ID NOs. 812-1599 comprising the steps of:

- a) reading said sequence through the use of a computer program which identifies features in sequences; and
- b) identifying features in said sequence with said computer program.

35

20. A vector comprising a nucleic acid according to either Claims 1 or 2.

21. A host cell containing a nucleic acid of Claim 20.

40

1/10

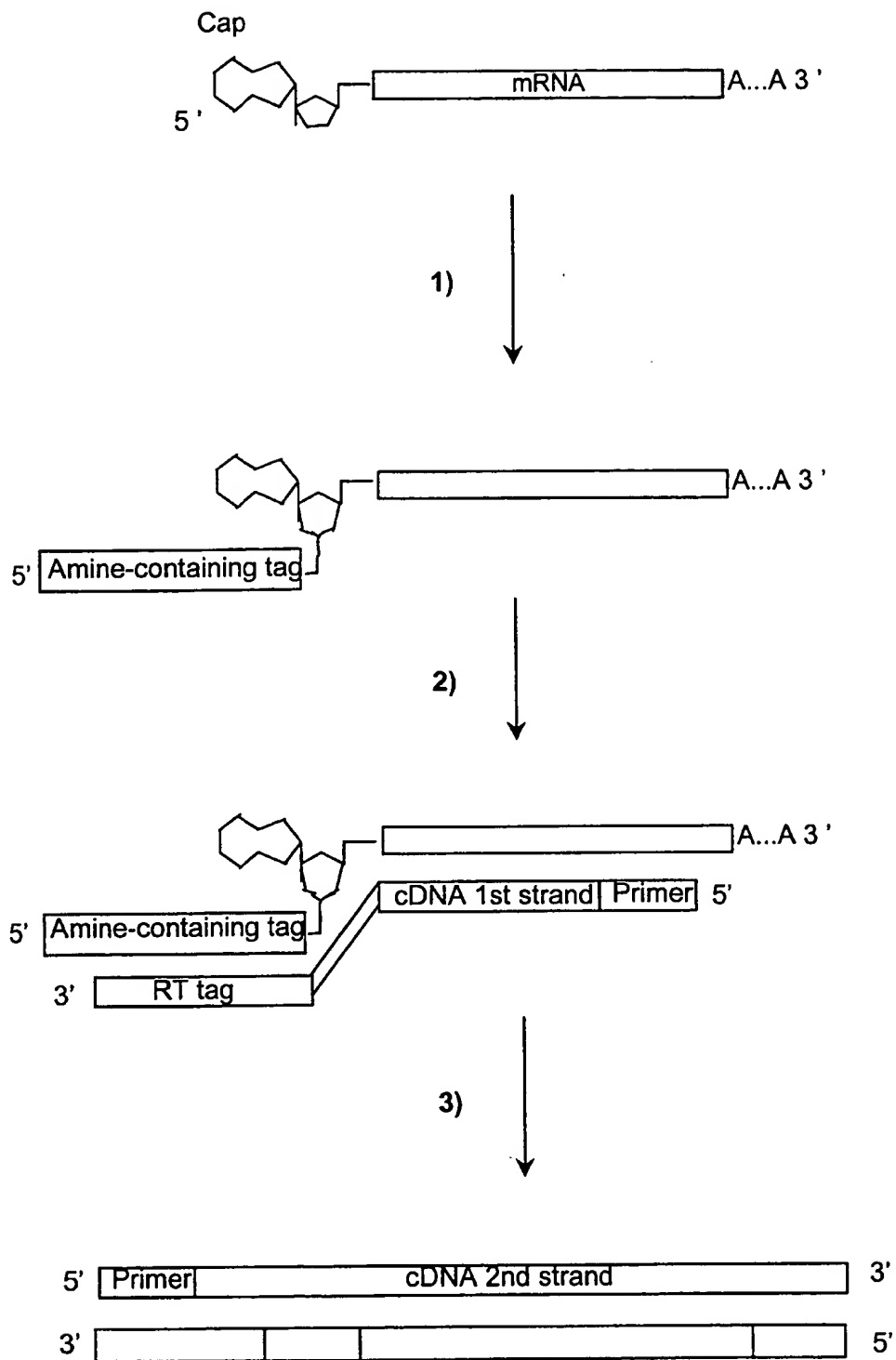


Figure 1

2/10

Minimum signal peptide score	false positive rate	false negative rate	proba(0.1)	proba(0.2)
3,5	0,121	0,036	0,467	0,664
4	0,096	0,06	0,519	0,708
4,5	0,078	0,079	0,565	0,745
5	0,062	0,098	0,615	0,782
5,5	0,05	0,127	0,659	0,813
6	0,04	0,163	0,694	0,836
6,5	0,033	0,202	0,725	0,855
7	0,025	0,248	0,763	0,878
7,5	0,021	0,304	0,78	0,889
8	0,015	0,368	0,816	0,909
8,5	0,012	0,418	0,836	0,92
9	0,009	0,512	0,856	0,93
9,5	0,007	0,581	0,863	0,934
10	0,006	0,679	0,835	0,919

Figure 2



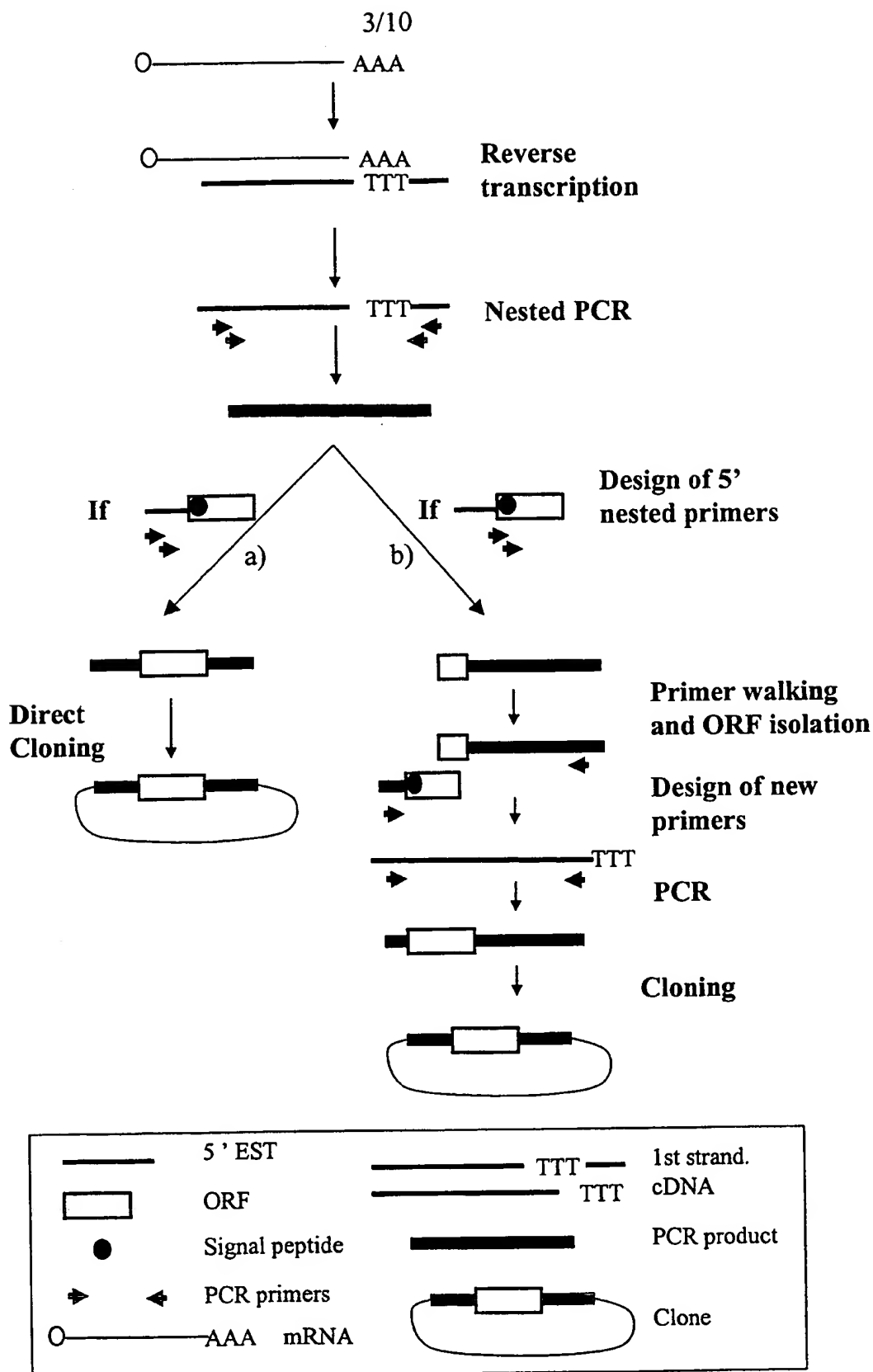


Figure 3

4/10

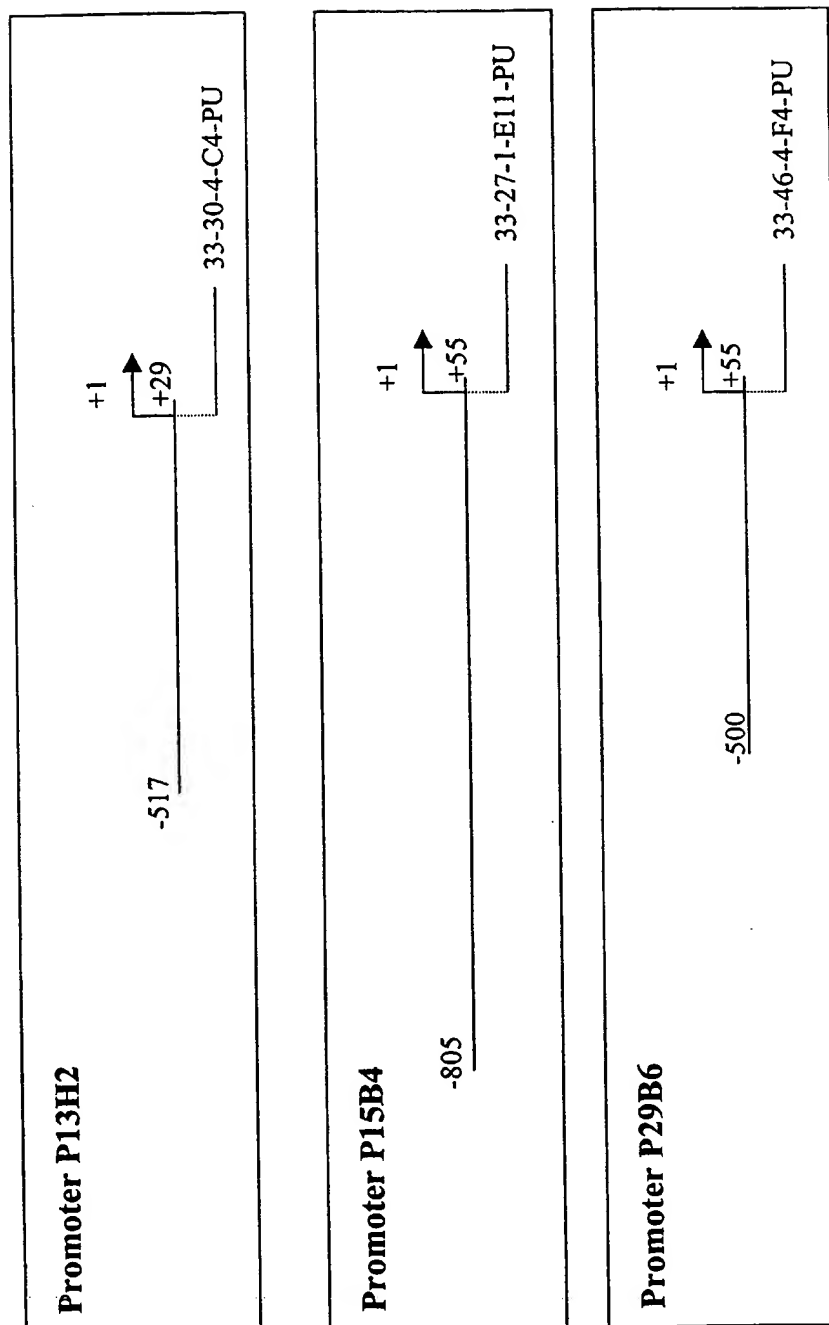


Figure 4

5/10

## Promoter sequence P13H2 (546 bp):

Matrix	Orient		Score	Length	Sequence
	Position	ation			
CMYB_01	-502	+	0.983	9	TGTCAGTTG
MYOD_Q6	-501	-	0.961	10	CCCAACTGAC
S8_01	-444	-	0.960	11	AATAGAATTAG
S8_01	-425	+	0.966	11	AACTAAATTAG
DELTAEF1_01	-390	-	0.960	11	GCACACCTCAG
GATA_C	-364	-	0.964	11	AGATAAATCCA
CMYB_01	-349	+	0.958	9	CTTCAGTTG
GATA1_02	-343	+	0.959	14	TTGTAGATAGGACA
GATA_C	-339	+	0.953	11	AGATAGGACAT
TAL1ALPHA47_01	-235	+	0.973	16	CATAACAGATGGTAAG
TAL1BETA47_01	-235	+	0.983	16	CATAACAGATGGTAAG
TAL1BETA1TF2_01	-235	+	0.978	16	CATAACAGATGGTAAG
MYOD_Q6	-232	-	0.954	10	ACCATCTGTT
GATA1_04	-217	-	0.953	13	TCAAGATAAAGTA
IK1_01	-126	+	0.963	13	AGTTGGGAATTCC
IK2_01	-126	+	0.985	12	AGTTGGGAATTC
CREL_01	-123	+	0.962	10	TGGGAATTCC
GATA1_02	-96	+	0.950	14	TCAGTGATATGGCA
SRY_02	-41	-	0.951	12	TAAACAAAACA
E2F_02	-33	+	0.957	8	TTAGCGC
MZF1_01	-5	-	0.975	8	TGAGGGGA

## Promoter sequence P15B4 (861bp) :

Matrix	Orient		Score	Length	Sequence
	Position	ation			
NFY_Q6	-748	-	0.956	11	GGACCAATCAT
MZF1_01	-738	+	0.962	8	CCTGGGGA
CMYB_01	-684	+	0.994	9	TGACCGTTG
VMYB_02	-682	-	0.985	9	TCCAACGGT
STAT_01	-673	+	0.968	9	TTCTTGGA
STAT_01	-673	-	0.951	9	TTCCAGGAA
MZF1_01	-556	-	0.956	8	TTGGGGGA
IK2_01	-451	+	0.965	12	GAATGGGATTC
MZF1_01	-424	+	0.986	8	AGAGGGGA
SRY_02	-398	-	0.955	12	GAAAACAAAACA
MZF1_01	-216	+	0.960	8	GAAGGGGA
MYOD_Q6	-190	+	0.981	10	AGCATCTGCC
DELTAEF1_01	-176	+	0.958	11	TCCCACCTTCC
S8_01	5	-	0.992	11	GAGGCAATTAT
MZF1_01	16	-	0.986	8	AGAGGGGA

## Promoter sequence P29B6 (555 bp) :

Matrix	Orient		Score	Length	Sequence
	Position	ation			
ARNT_01	-311	+	0.964	16	GGACTCACGTGCTGCT
NMYC_01	-309	+	0.965	12	ACTCACGTGCTG
USF_01	-309	+	0.985	12	ACTCACGTGCTG
USF_01	-309	-	0.985	12	CAGCACGTGAGT
NMYC_01	-309	-	0.956	12	CAGCACGTGAGT
MYCMAX_02	-309	-	0.972	12	CAGCACGTGAGT
USF_C	-307	+	0.997	8	TCACGTGC
USF_C	-307	-	0.991	8	GCACGTGA
MZF1_01	-292	-	0.968	8	CATGGGGA
ELK1_02	-105	+	0.963	14	CTCTCCGGAAGCCT
CETS1P54_01	-102	+	0.974	10	TCCGGAAGCC
AP1_Q4	-42	-	0.963	11	AGTGACTGAAC
AP1FJ_Q2	-42	-	0.961	11	AGTGACTGAAC
PADS_C	45	+	1.000	9	TGTGGTCTC

Figure 5

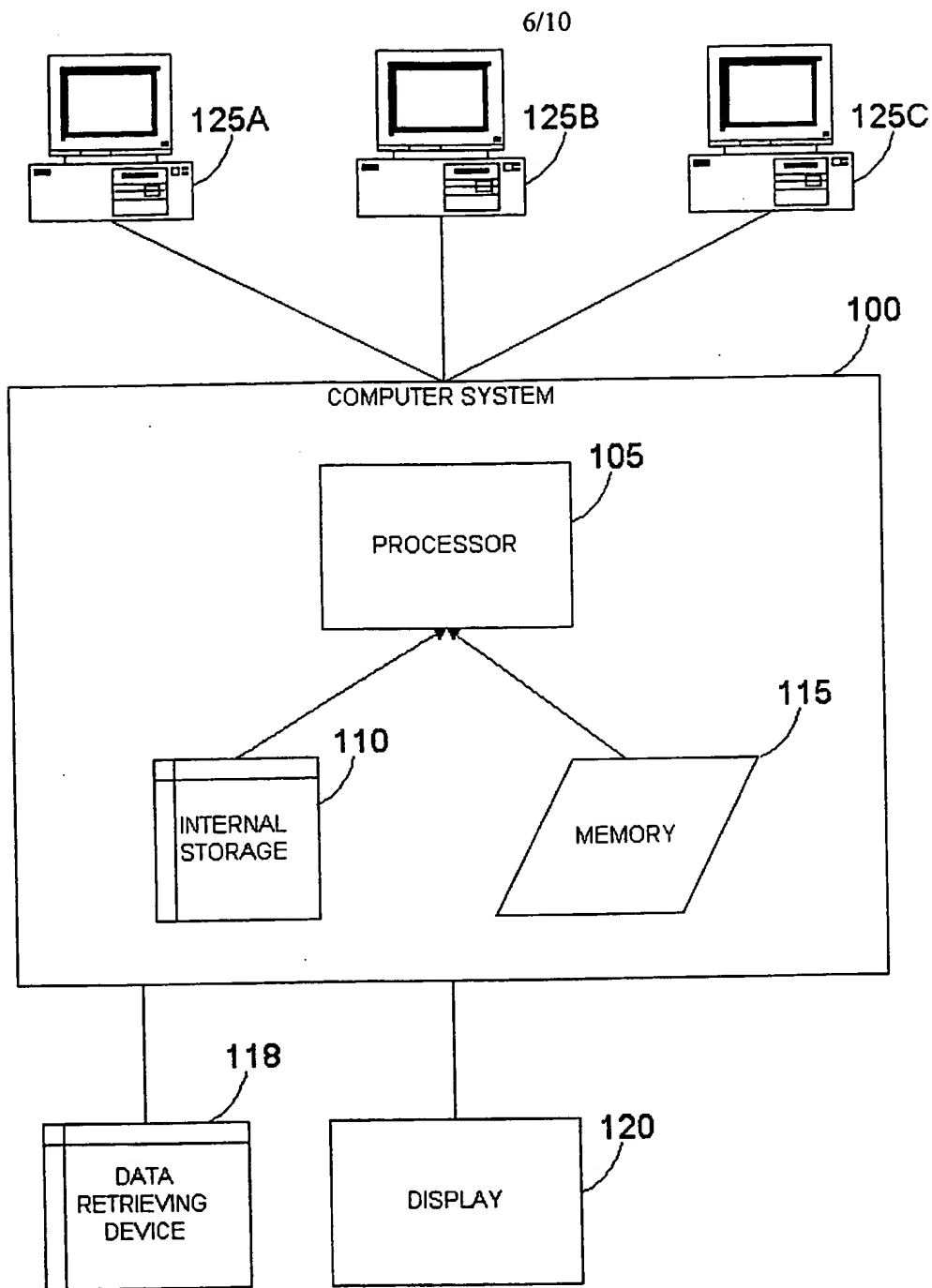


FIGURE 6

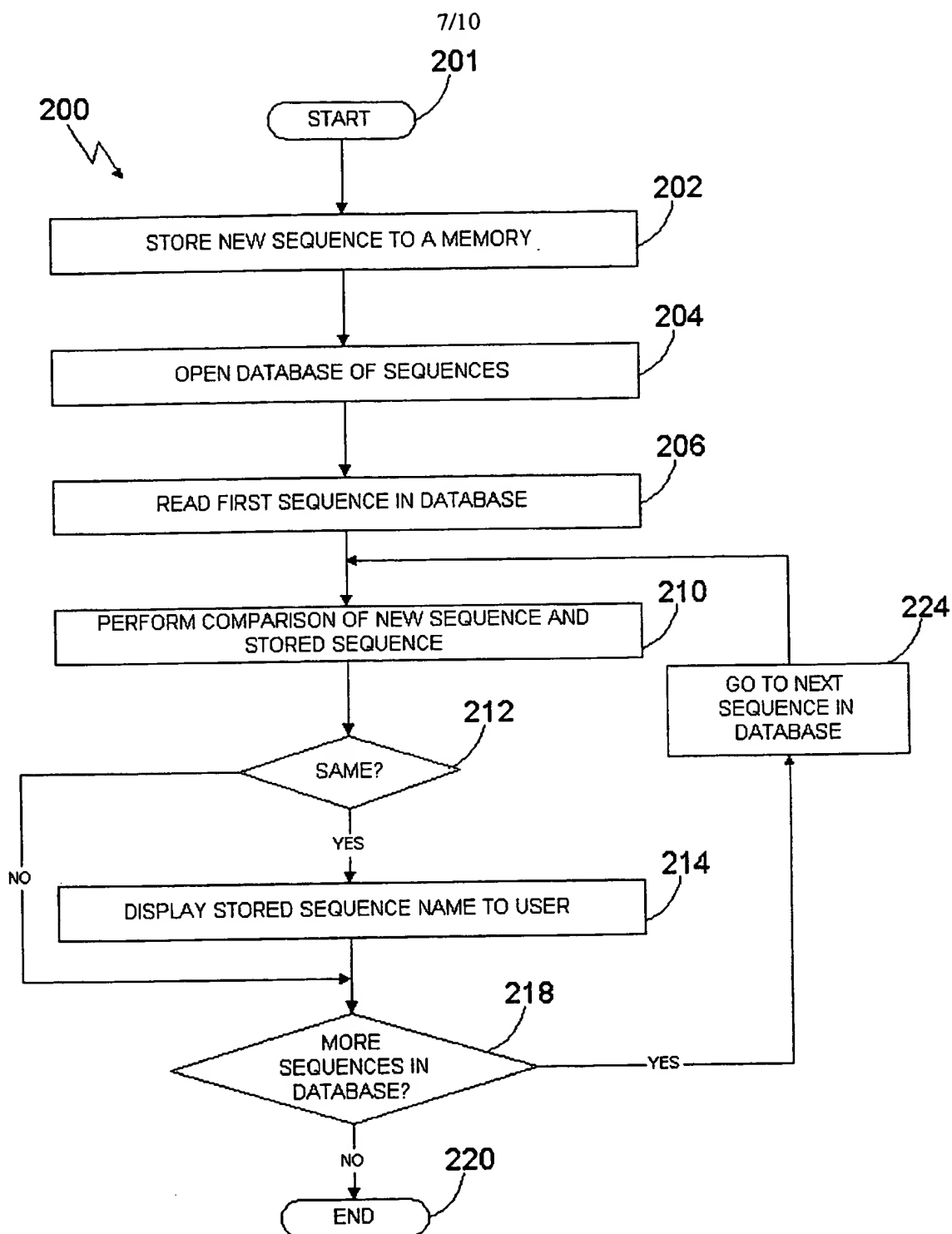


FIGURE 7

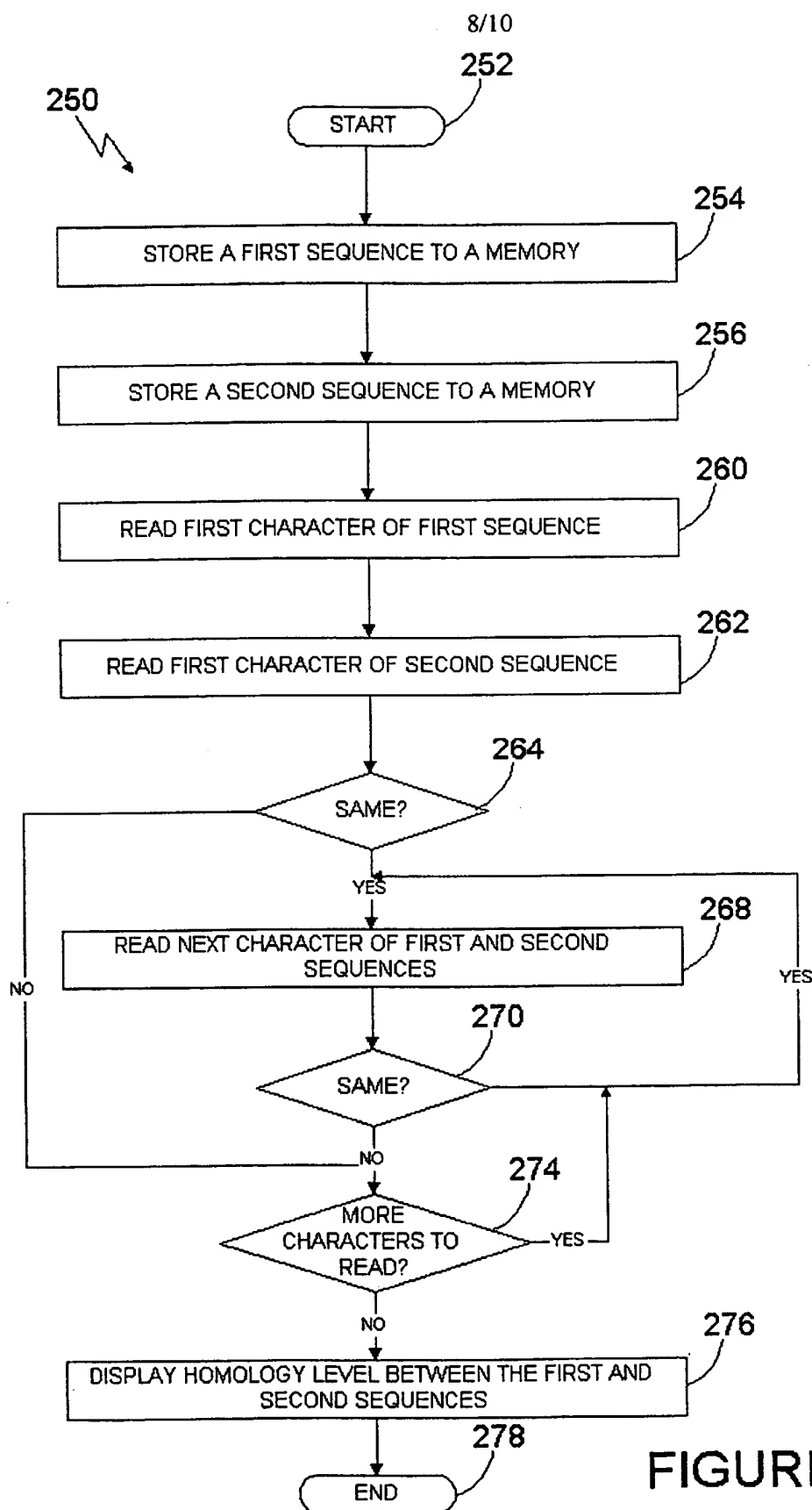


FIGURE 8

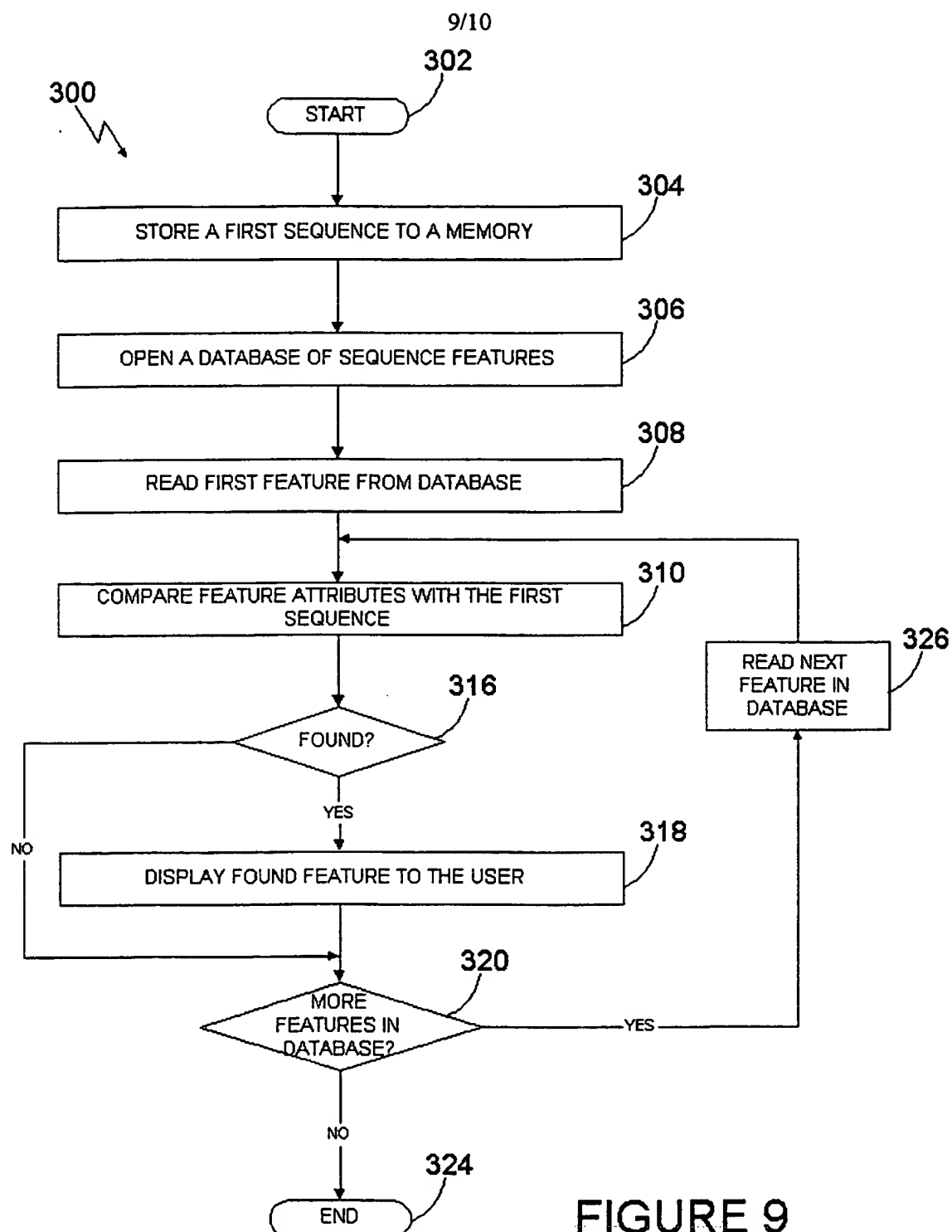


FIGURE 9

Step	Search characteristic		Selection Characteristics			
	Program	Strand	Parameters	Identity (%)	Length (bp)	Comments
miscellaneous	FASTA	both	-	90	15	
tRNA	FASTA	both	-	80	60	
rRNA	BLASTN	both	S=108	80	40	
mtRNA	BLASTN	both	S=108	80	40	
Procaryotic	BLASTN	both	S=144	90	40	
Fungal	BLASTN	both	S=144	90	40	
Alu	BLASTN	both	S=72, B=5	70	40	max 5 matches, masking
L1	BLASTN	both	S=72, B=5	70	40	max 5 matches, masking
Repeats	BLASTN	both	S=72	70	40	masking
PolyA	BLAST2N	top	W=6, S=10, E=1000, N=12	90	10	in the last 100 nucleotides
Polyadenylation signal	-	top	AATAAA allowing 1 mismatch	90 then 70	30	in the 50 nucleotides before the 5' end of the polyA
Vertebrate	BLASTN then FASTA	both	-	90 then 70	30	first BLASTN, then FASTA on matching sequences
ESTs	BLAST2N	both	-	90	30	
Geneseq	BLASTN	both	W=8, B=10	90	30	
ORF	BLASTP	top	W=8, B=10	-	-	on ORF proteins, max 10 matches
Proteins	BLASTX	top	E = 0.001	70	30	

Figure 10



## SEQUENCE LISTING

<110> Dumas Milne Edwards, J.B.  
 Duclert A.  
 Giordano, J.Y.  
 Genset SA

<120> ESTs and Encoded Human Proteins.

<130> D18118-339 881

<150> 09/057,719

<151> 1998-04-09

<150> 09/069,047

<151> 1998-04-28

<160> 1622

<170> Patent.pm

<210> 1

<211> 822

<212> DNA

<213> Homo Sapiens

<220>

<221> CDS

<222> 346..552

<221> sig\_peptide

<222> 346..408

<223> Von Heijne matrix

<221> misc\_feature

<222> 115

<223> n=a, g, c or t

<400> 1

```

actcctttta gcataggggc ttccggcgcca gcggccagcg ctagtcgggc tggtaagtgc      60
ctgatgccga gtccgctctc tcgctgtctt tcttggtccc aggcaaagcg gasgnagatc      120
ctcaaacggc ctagtgcttc gcgcttccgg agaaaatcag cgggtctaatt aattcctctg      180
gtttgttgaa gcagttacca agaattctca accctttccc acaaaagcta attgagtaca      240
cgttcctggt gagtacacgt tctgttgat ttacaaaagg tgcaggatg agcaggctctg      300
aagactaaca ttttgtgaag ttgtaaaaca gaaaacctgt tagaa atg tgg tgg ttt      357
                                     Met Trp Trp Phe
                                     -20
cag caa ggc ctc agt ttc ctt cct tca gcc ctt gta att tgg aca tct      405
Gln Gln Gly Leu Ser Phe Leu Pro Ser Ala Leu Val Ile Trp Thr Ser
      -15                               -10                               -5
gct gct ttc ata ttt tca tac att act gca gta aca ctc cac cat ata      453
Ala Ala Phe Ile Phe Ser Tyr Ile Thr Ala Val Thr Leu His His Ile
      1                               5                               10                               15
gac ccg gct tta cct tat atc agt gac act ggt aca gta gct cca raa      501
Asp Pro Ala Leu Pro Tyr Ile Ser Asp Thr Gly Thr Val Ala Pro Xaa
      20                               25                               30
aaa tgc tta ttt ggg gca atg cta aat att gcg gca gtt tta tgt caa      549
Lys Cys Leu Phe Gly Ala Met Leu Asn Ile Ala Ala Val Leu Cys Gln
      35                               40                               45
aaa tagaaatcag gaarataatt caacttaaag aakttcattt catgacaaa      602
Lys
ctcttcaraa acatgtcttt acaagcatat ctcttgattt gctttctaca ctgttgaatt      662
gtctggcaat atttctgcag tggaaaattt gatttarmta gttcttgact gataaatatg      722

```

&lt;222&gt; 25..84

&lt;223&gt; Von Heijne matrix

score 11.1000003814697

seq LLALLFFLGQAAG/DL

&lt;400&gt; 46

```

agcggctcca gctaagagga caag atg agg ccc ggc ctc tca ttt ctc cta      51
                        Met Arg Pro Gly Leu Ser Phe Leu Leu
                        -20                      -15

gcc ctt ctg ttc ttc ctt ggc caa gct gca ggg gat ttg ggg gat gtg      99
Ala Leu Leu Phe Phe Leu Gly Gln Ala Ala Gly Asp Leu Gly Asp Val
-10                      -5                      1                      5
gga cct cca att ccc agc ccc ggc ttc agc tct ttc cca ggt gtt gac      147
Gly Pro Pro Ile Pro Ser Pro Gly Phe Ser Ser Phe Pro Gly Val Asp
                        10                      15                      20
tcc agc tcc agc ttc agc tcc agc tcc agg tcg ggc tcc agc tcc agc      195
Ser Ser Ser Ser Phe Ser Ser Ser Ser Arg Ser Gly Ser Ser Ser Ser
                        25                      30                      35
cgc agc tta ggc agc gga ggt tct gtg tcc cag ttg ttt tcc aat ttc      243
Arg Ser Leu Gly Ser Gly Gly Ser Val Ser Gln Leu Phe Ser Asn Phe
                        40                      45                      50
acc ggc tcc gtg gat gac cgt ggg acc tgc cag tgc tct gtt tcc ctg      291
Thr Gly Ser Val Asp Asp Arg Gly Thr Cys Gln Cys Ser Val Ser Leu
                        55                      60                      65
cca gac acc acc ttt ccc gtg gac aga gtg gaa cgc ttg gaa ttc aca      339
Pro Asp Thr Thr Phe Pro Val Asp Arg Val Glu Arg Leu Glu Phe Thr
70                      75                      80                      85
gct cat gtt ctt tct cag aag ttt gag aaa gaa ctt tct aaa gc      383
Ala His Val Leu Ser Gln Lys Phe Glu Lys Glu Leu Ser Lys
                        90                      95

```

&lt;210&gt; 47

&lt;211&gt; 459

&lt;212&gt; DNA

&lt;213&gt; Homo sapiens

&lt;220&gt;

&lt;221&gt; CDS

&lt;222&gt; 17..457

&lt;221&gt; sig\_peptide

&lt;222&gt; 17..94

&lt;223&gt; Von Heijne matrix

score 11.1000003814697

seq FLLLVAAPRWVRS/QV

&lt;221&gt; misc\_feature

&lt;222&gt; 399

&lt;223&gt; n=a, g, c or t

&lt;400&gt; 47

```

atactttctg agactc atg gac ctc ctg cac aag aac atg aaa cac ctg tgg      52
                        Met Asp Leu Leu His Lys Asn Met Lys His Leu Trp
                        -25                      -20                      -15

ttc ttc ctc ctc ctg gtg gca gct ccc aga tgg gtc cgg tct car gtg      100
Phe Phe Leu Leu Leu Val Ala Ala Pro Arg Trp Val Arg Ser Gln Val
-10                      -5                      1
cag ctg cak gag tcg ggc cca gga ctg gtg aag cct tcg ggg acc ctg      148
Gln Leu Xaa Glu Ser Gly Pro Gly Leu Val Lys Pro Ser Gly Thr Leu
5                      10                      15
tcc ctc atc tgc ggt gtc tct ggt gat tcc gtc acc att agt ggt tgg      196
Ser Leu Ile Cys Gly Val Ser Gly Asp Ser Val Thr Ile Ser Gly Trp
20                      25                      30

```

34

```

tgg agt tgg gtc cgc cag ccc cca ggg aag gga ctg gag tgg att tcg      244
Trp Ser Trp Val Arg Gln Pro Pro Gly Lys Gly Leu Glu Trp Ile Ser
35          40          45          50
gaa atc gat cat ggt gga aac acc aat tac aac ccg tcc ctc aag agt      292
Glu Ile Asp His Gly Gly Asn Thr Asn Tyr Asn Pro Ser Leu Lys Ser
          55          60          65
cga gtc kcc att tct tta gac aag tcc aag aat aag ttc tcc ctg agg      340
Arg Val Xaa Ile Ser Leu Asp Lys Ser Lys Asn Lys Phe Ser Leu Arg
          70          75          80
ctg acc tct gtg acc gcc gcg gac acc gcc atg tat kac tgt gcg aga      388
Leu Thr Ser Val Thr Ala Ala Asp Thr Ala Met Tyr Xaa Cys Ala Arg
          85          90          95
ggc ggt gcg bnc agc tcg tcc gct ttt gat gtc tgg ggc cta rgg aca      436
Gly Gly Ala Xaa Ser Ser Ser Ala Phe Asp Val Trp Gly Leu Xaa Thr
          100          105          110
atg gtc atc atc tct tca gcc tc      459
Met Val Ile Ile Ser Ser Ala
115          120

```

&lt;210&gt; 48

&lt;211&gt; 437

&lt;212&gt; DNA

&lt;213&gt; Homo sapiens

&lt;220&gt;

&lt;221&gt; CDS

&lt;222&gt; 20..436

&lt;221&gt; sig\_peptide

&lt;222&gt; 20..76

&lt;223&gt; Von Heijne matrix

score 11

seq TLLLLLTVPSWVLS/QV

&lt;400&gt; 48

```

gtgaatcctg ctctccacc atg gac ata ctt tgt tcc acg ctc ctg ctm ctg      52
Met Asp Ile Leu Cys Ser Thr Leu Leu Leu Leu
          -15          -10
ack gtc ccg tcc tgg gtc tta tcc car gtc acc ttg arg gaa tct ggt      100
Thr Val Pro Ser Trp Val Leu Ser Gln Val Thr Leu Xaa Glu Ser Gly
          -5          1          5
cct gcg ctg gtg aaa gcc aca cag acc ctc aga ctg acc tgc acc ttc      148
Pro Ala Leu Val Lys Ala Thr Gln Thr Leu Arg Leu Thr Cys Thr Phe
          10          15          20
tct ggg ttc tca ctc agc act aat aga atg cgt gtg agt tgg atc cgt      196
Ser Gly Phe Ser Leu Ser Thr Asn Arg Met Arg Val Ser Trp Ile Arg
25          30          35          40
cag ccc cca ggg aag gcc ctg gag tgg ctt gca cgg att gat tgg gat      244
Gln Pro Pro Gly Lys Ala Leu Glu Trp Leu Ala Arg Ile Asp Trp Asp
          45          50          55
gat tat aag agg tac agc aca tct ctg aag acc agg gtc acc atc tcc      292
Asp Tyr Lys Arg Tyr Ser Thr Ser Leu Lys Thr Arg Val Thr Ile Ser
          60          65          70
aag gac acg tcc aaa aac cag gtg atc ctg aca atg acc aac gtg gac      340
Lys Asp Thr Ser Lys Asn Gln Val Ile Leu Thr Met Thr Asn Val Asp
          75          80          85
cct gcg gac aca gcc acc tat tac tgt gca cgc ctt tca acg gca gct      388
Pro Ala Asp Thr Ala Thr Tyr Tyr Cys Ala Arg Leu Ser Thr Ala Ala
          90          95          100
acc cca cag ttt ttt gac ttc tgg ggc cag gga gtc ctg gtc tcc gtc t      437
Thr Pro Gln Phe Phe Phe Thr Trp Gly Gln Gly Val Leu Val Ser Val
105          110          115          120

```

<221> SIGNAL  
 <222> -20...-1

<400> 833

```

Met Glu Lys Ile Pro Val Ser Ala Phe Leu Leu Leu Val Ala Leu Ser
-20          -15          -10          -5
Tyr Thr Leu Ala Arg Asp Thr Thr Val Lys Pro Gly Ala Lys Lys Asp
          1          5          10
Thr Lys Asp Ser Arg Pro Lys Leu Pro Gln Thr Leu Ser Arg Gly Trp
          15          20          25
Gly Asp Gln Leu Ile Trp Thr Gln Thr Tyr Glu Glu Ala Leu Tyr Lys
          30          35          40
Ser Lys Thr Ser Asn Lys Pro Leu Met Ile Ile His His Leu Asp Glu
45          50          55          60
Cys Pro His Ser Gln Ala Leu Lys Lys Val Phe Ala Glu Asn Lys Glu
          65          70          75
Ile Gln Lys Leu Ala Glu Gln Phe Val Leu Leu Asn Leu Val Tyr Glu
          80          85          90
Thr Thr Asp
          95
  
```

<210> 834  
 <211> 119  
 <212> PRT  
 <213> Homo sapiens

<220>  
 <221> SIGNAL  
 <222> -20...-1

<400> 834

```

Met Arg Pro Gly Leu Ser Phe Leu Leu Ala Leu Leu Phe Phe Leu Gly
-20          -15          -10          -5
Gln Ala Ala Gly Asp Leu Gly Asp Val Gly Pro Pro Ile Pro Ser Pro
          1          5          10
Gly Phe Ser Ser Phe Pro Gly Val Asp Ser Ser Ser Ser Phe Ser Ser
          15          20          25
Ser Ser Arg Ser Gly Ser Ser Ser Arg Ser Leu Gly Ser Gly Gly
          30          35          40
Ser Val Ser Gln Leu Phe Ser Asn Phe Thr Gly Ser Val Asp Asp Arg
45          50          55          60
Gly Thr Cys Gln Cys Ser Val Ser Leu Pro Asp Thr Thr Phe Pro Val
          65          70          75
Asp Arg Val Glu Arg Leu Glu Phe Thr Ala His Val Leu Ser Gln Lys
          80          85          90
Phe Glu Lys Glu Leu Ser Lys
          95
  
```

<210> 835  
 <211> 147  
 <212> PRT  
 <213> Homo sapiens

<220>  
 <221> SIGNAL  
 <222> -26...-1

<400> 835

```

Met Asp Leu Leu His Lys Asn Met Lys His Leu Trp Phe Phe Leu Leu
-25          -20          -15
Leu Val Ala Ala Pro Arg Trp Val Arg Ser Gln Val Gln Leu Xaa Glu
-10          -5          1          5
Ser Gly Pro Gly Leu Val Lys Pro Ser Gly Thr Leu Ser Leu Ile Cys
  
```

Identifier: AAZ42287 cDNA Sequence 383 BP  
Release Info: Derwent Geneseq Database Release No. 200202; Date released 21-JAN-02  
Database XReference: WPI; 2000-038446/03.;P-PSDB; AAY64673.  
Accession Number: AAZ42287  
Patent Title: Novel secreted protein 5' expressed sequence tag sequences used in diagnostic, forensic, gene therapy, and chr  
Patented by: (GEST ) GENSET.  
Inventor: Dumas Milne Edwards J, Duclert A, Giordano J  
Description: Human 5' EST isolated from a cDNA library SEQ ID NO:46.  
Patent Number: WO9953051-A2  
Patent Publication Date: 21-OCT-1999  
Modification Date: 01-FEB-2000 (first entry)  
Local Filing: 09-APR-1999; 99WO-IB00712  
Priority: 09-APR-1998  
Abstract: AAZ42265 to AAZ43075 represent novel 5' expressed sequence tag (EST) sequences, corresponding to huma untranslated regions (UTRs) and upstream regulatory regions which control the location, development stage, r in diagnostic procedures to identify individuals having genetic diseases resulting from abnormal gene expressi polypeptide into a cell. The proteins encoded by the EST sequences may be useful in treating a variety of hum  
KeyWords: Human;5' EST;expressed sequence tag;secreted protein;diagnosis;gene therapy;chromosome mapping;upstreai  
Organism: Homo sapiens.  
Sequence Composition: Sequence 383 BP; 69 A; 118 C; 98 G; 98 T; 0 other;  
Sequence: >AAZ42287 WO9953051-A2 PA (GEST ) PR 09-APR-1998 PF 09-APR-1999 Human 5' EST isolated from :  
TTCTTCCTTGGCCAAGCTGCAGGGGATTTGGGGGATGTGGGACCTCCAATTCCCAGCCCC GGCTT  
TTCACGGGCTCCGTGGATGACCGTGGGACCTGCCAGTGCTCTGTTTCCCTGCCAGACACC ACCTT

Identifier: AAY64673 Protein Sequence 119 AA  
 Release Info: Derwent Geneseq Database Release No. 200202; Date released 21-JAN-02  
 Database XReference: WPI; 2000-038446/03.;N-PSDB; AAZ42287.  
 Accession Number: AAY64673  
 Patent Title: Novel secreted protein 5' expressed sequence tag sequences used in diagnostic, forensic, gene therapy, and chromosome mapping procedures  
 Patented by: (GEST ) GENSET.  
 Inventor: Dumas Milne Edwards J, Duclert A, Giordano J  
 Description: Human 5' EST related polypeptide SEQ ID NO:834.  
 Patent Number: WO9953051-A2  
 Patent Publication Date: 21-OCT-1999  
 Modification Date: 01-FEB-2000 (first entry)  
 Local Filing: 09-APR-1999; 99WO-IB00712  
 Priority: 09-APR-1998  
 Abstract: AAZ42265 to AAZ43075 represent novel 5' expressed sequence tag (EST) sequences, corresponding to human secreted proteins. AAY64651 to AAY65438 represent the EST-related proteins corresponding to AAZ42265 to AAZ43075. The 5' ESTs can be used for producing secreted human gene products. They can be used to identify and isolate 5' untranslated regions (UTRs) and upstream regulatory regions which control the location, development, stage, rate, and quantity of protein synthesis, as well as stability of mRNA. The ESTs are also useful as probes for chromosome mapping, and to obtain full length cDNA clones. The ESTs can also be used in forensic procedures to identify individuals, or in diagnostic procedures to identify individuals having genetic diseases resulting from abnormal gene expression. The products may also be used in gene therapy protocols. The nucleic acids encoding signal peptides can be used for directing extracellular secretion of a polypeptide or the insertion of a polypeptide into a membrane, or importing a polypeptide into a cell. The proteins encoded by the EST sequences may be useful in treating a variety of human conditions. Secreted proteins have therapeutic value, and the identification of new secreted proteins is valuable. AAZ42249 to AAZ42264 and AAY64644 to AAY64650 represent sequences used as the exemplification of the present invention.  
 KeyWords: Human;5' EST;expressed sequence tag;secreted protein;diagnosis;gene therapy;chromosome mapping;upstream regulatory sequence;forensic;location;development;protein synthesis;stability;regulation;identification.  
 Organism: Homo sapiens.  
 Sequence Composition: Sequence 119 AA; 4 A; 6 R; 1 N; 7 D; 0 B; 2 C; 4 Q; 4 E; 0 Z; 13 G; 1 H; 1 I; 13 L; 3 K; 1 M; 11 F; 8 P; 27 S; 0 Y; 8 V; 0 Others;  
 Sequence: >AAY64673 WO9953051-A2 PA (GEST ) PR 09-APR-1998 PF 09-APR-1999 Human 5' EST related polypeptide SEQ ID NO:834. [Homo sapiens.]  
 MRPGLSFLALLFFLGQAAGDLGDVGPPISPFGSSFPGVDSFFFSSSSRSRSSSL